

Machine Learning in Large Scale Structure

Shirley Ho (Carnegie Mellon University)

Work with Xiaoying Xu,
Andrea Klein (CS), Shadab Alam, Zongge Liu
Jeff Schneider (CS), Barnabas Poczos (CS), Junier Oliver (CS)

Carnegie Mellon University

What do you see on campus?

Carnegie Mellon University What do you see on campus?



Bat vehicle + rioting Gotham citizens ...

Carnegie Mellon University

What do you see on campus?



Cat lady escaping from
the Dean's office

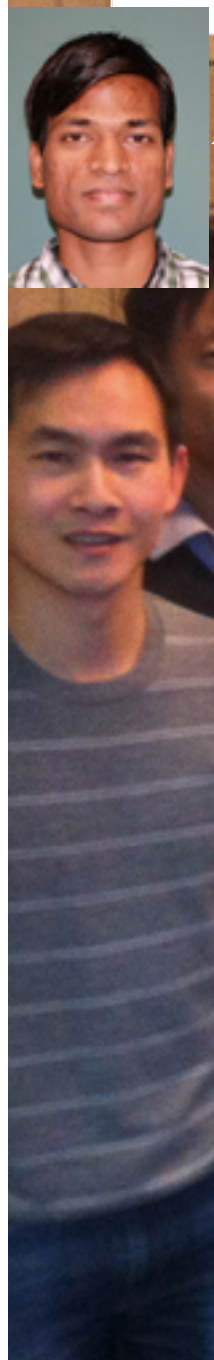


you also see nice happy cosmologist ...

Machine Learning

Dark Energy

Inflation



Outline

- Simulation Production
 - Populating Halos with Galaxies using Machine Learning
 - Generating the Density field

Observations:

flux(x,y,band/wavelength), observational systematics(x,y,band/wavelength): psf, sky, dust, airmass and respective errors

Basic Reduction pipeline

Extracted object, observational systematics properties

LSS systematics removal, statistics estimation

Large Scale Structure (BAO, full clustering measurements)

Theory, Beyond Linear
Model Predictions

Simulations (possibly with
DM simulations + HOD),
Covariance Matrix

Cosmology (cosmological parameters, formation and
evolution of galaxies, quasars...)

Observations:

flux(positions,band/wavelength), observational systematics(positions,band/wavelength):
psf, sky, dust, airmass and respective errors

Basic Reduction pipeline

Extracted object, observational systematics properties

LSS systematics removal, statistics estimation

Large Scale Structure measurements
BAO (with reconstruction), full clustering

Linear Theory and
Beyond Linear Modeling

Covariance Matrix

Astronomy and Cosmology
[cosmological parameters, formation and evolution of galaxies, quasars...]

Why Populate Halos?

- Mock galaxy catalogues! These are necessary for LSS analysis
 - calculating the covariance matrix, testing the pipeline, etc.
- Alternatives:
 - Run N-body + hydro simulations, very expensive if you want many mocks.
 - Perturbation theory, not accurate enough on scales $< 20\text{Mpc}$, especially in redshift space.

Why Populate Halos?

- Mock galaxy catalogues! These are necessary for LSS analysis
 - calculating the covariance matrix, testing the pipeline, etc.
- Alternatives:
 - Run N-body + hydro simulations, very expensive if you want many mocks.
 - Perturbation theory, not accurate enough on scales $< 20\text{Mpc}$, especially in redshift space.

NB:

Here, by populating halos, I specifically mean determining the number of galaxies that will reside in a halo given its properties.

Why Populate Halos?

- Mock galaxy catalogues! These are necessary for LSS analysis
 - calculating the covariance matrix, testing the pipeline, etc.
- Alternatives:
 - Run N-body + hydro simulations, very expensive if you want many mocks.
 - Perturbation theory, not accurate enough on scales $< 20\text{Mpc}$, especially in redshift space.

Introducing Machine Learning

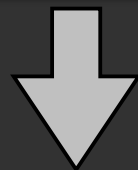
- Machine Learning algorithms **learns trends from the data itself**, it does not impose pre-assumed models on the data.
- The advantage of ML is that it is fully non-parametric:
 - The only assumption necessary is that some relationship **does** exist between halo properties (features) and the number of galaxies that will reside in it and that this relationship is continuous.

Cool Examples of ML applications (Courtesy Slide from Kayvon Fatahalian)

[Shrivastava 2011]

“Find images that are similar to a query image (even if not similar in individual pixel values).”

Query image
(snowy day)



Matches



[Doersch 2012]

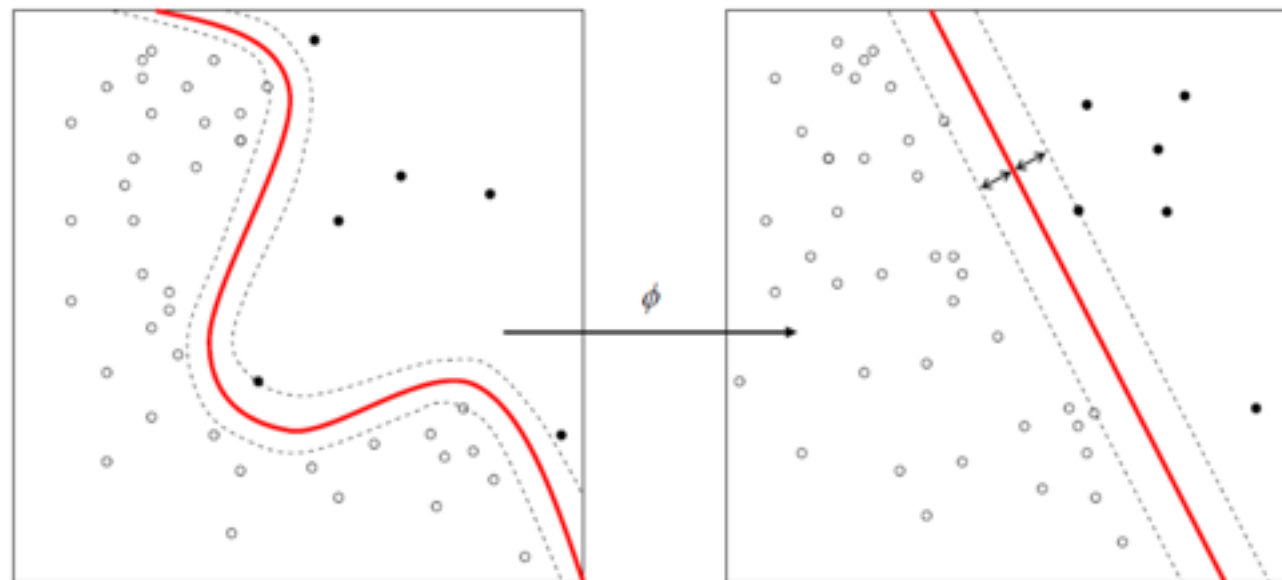
“Find meaningful visual elements that are unique to Paris”

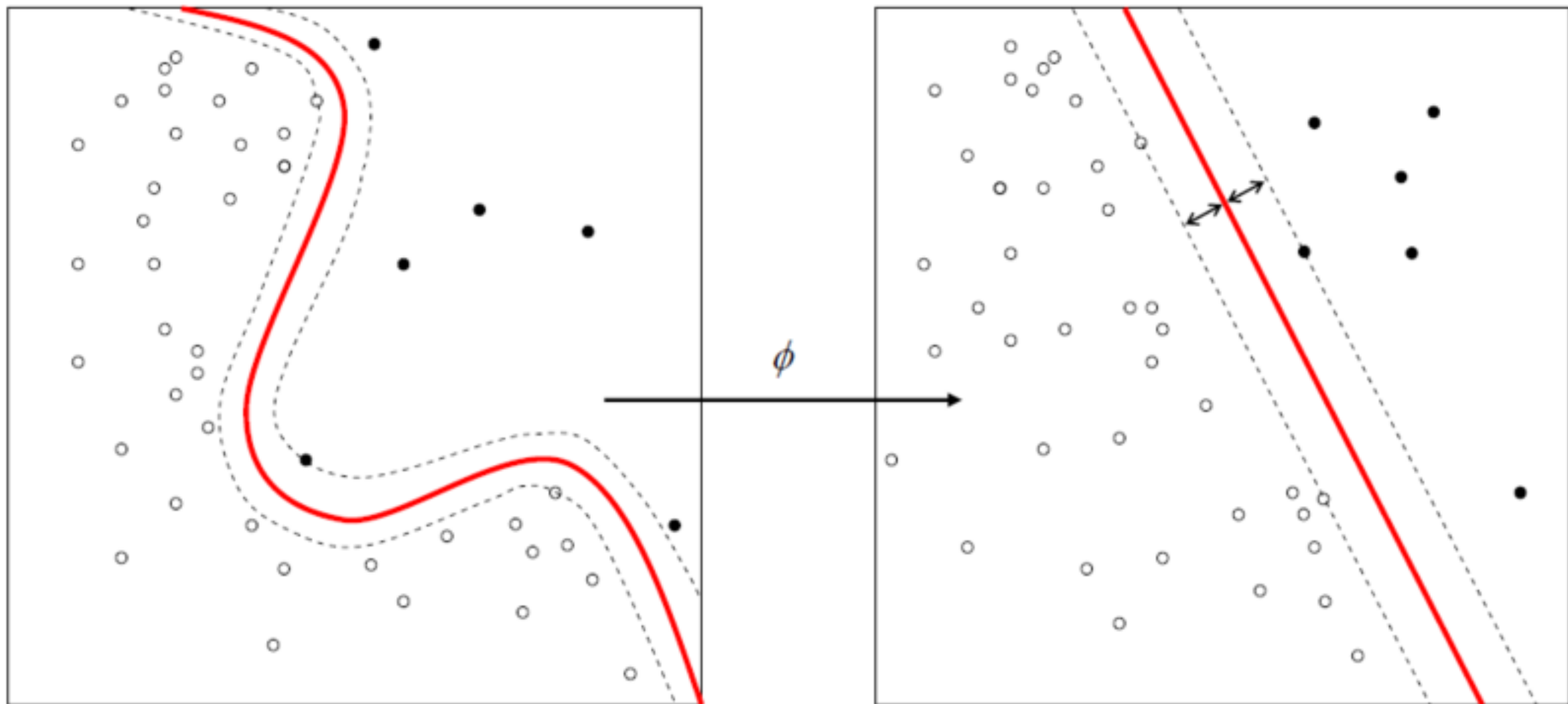


Simple ML algorithms

SVM: Support Vector Machine

- SVM maps the data into higher dimensional space with a kernel. To train, separate data into classes using hyperplanes. It can be generalized into regression (not only classification).





Simple ML algorithms

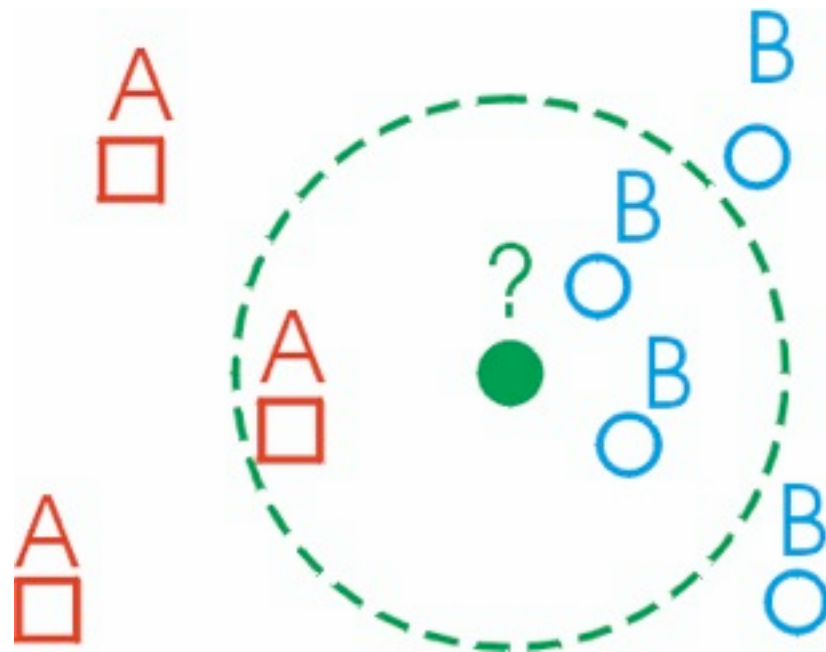
kNN: k-nearest neighbors

- kNN takes the average of k nearest neighbors to the point of interest in the training set.

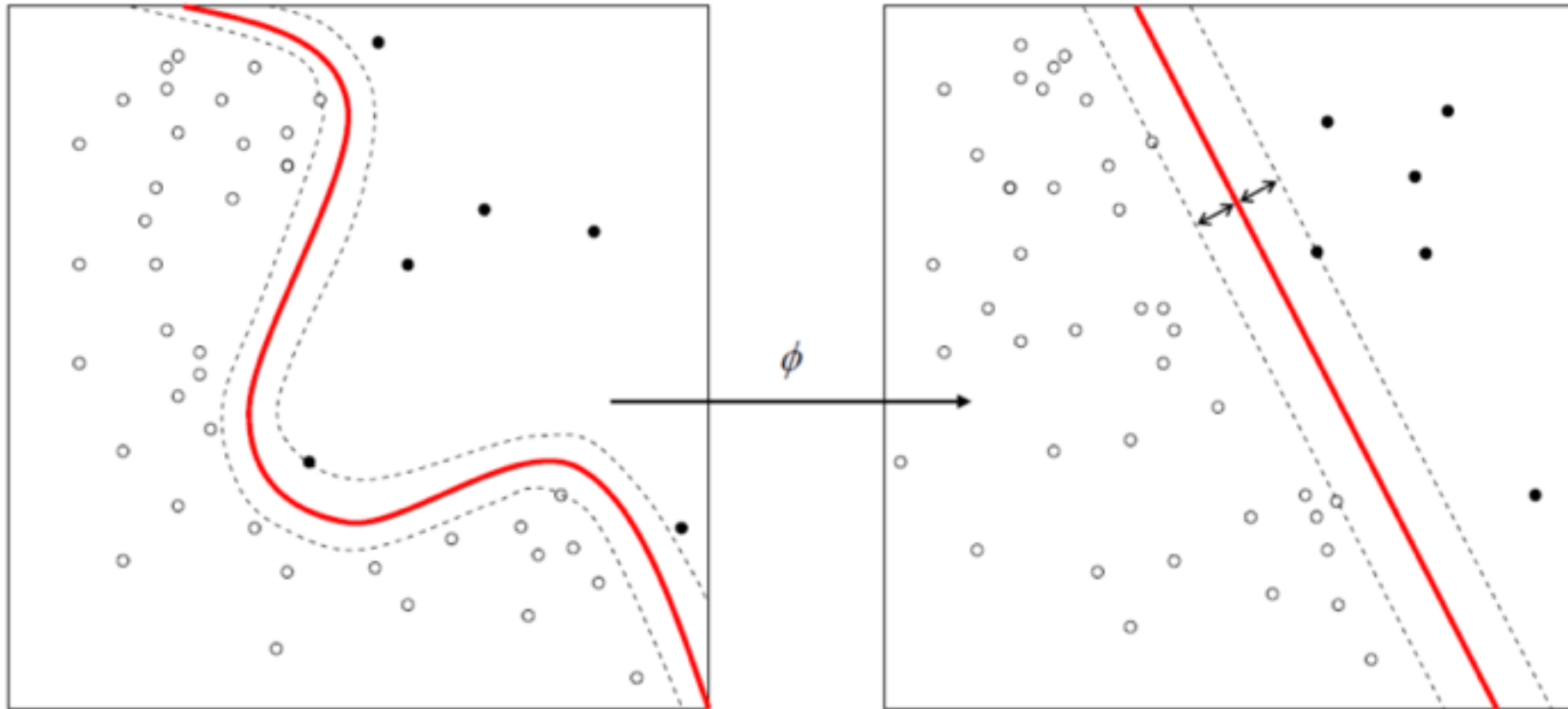
Simple ML algorithms

kNN: k-nearest neighbors

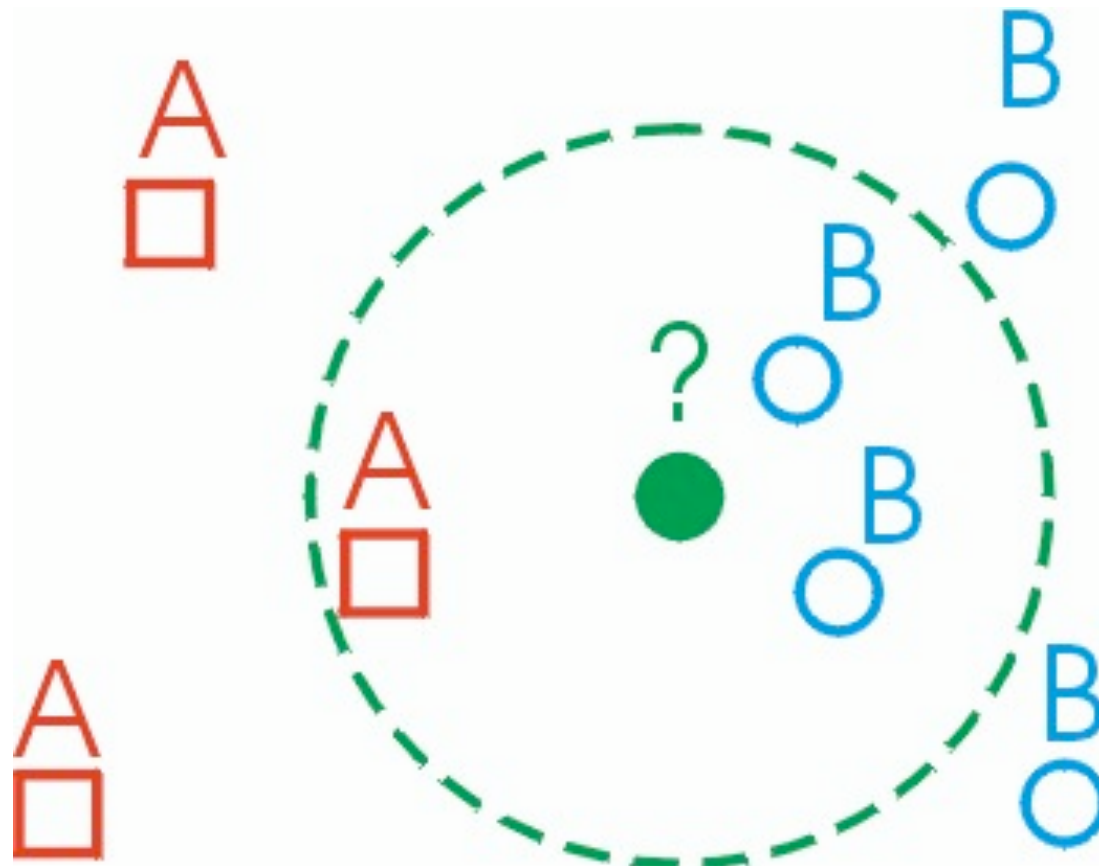
- kNN takes the average of k nearest neighbors to the point of interest in the training set.



SVM



kNN



Introducing Machine Learning (cont'd)

- To evaluate how good a machine learning algorithm perform, we use the “mean-squared-error (MSE):
- $$\text{MSE} = \frac{\sum_{i=1}^N (Y_{i,\text{test,true}} - Y_{i,\text{test,predicted}})^2}{N}$$
- We can compare the MSEs given by different algorithms and see which one does a better job, and we can also compare it to the **base MSE**, which basically replace $Y_{\text{test,predicted}}$ by average of Y_{train}

Machine Learning

Populating Halos with Galaxies

- We want to use Machine Learning to learn about how many galaxies (or specific kinds of galaxies) are in halos with certain properties.
- We can do this with real data, however, we are not in the era where we have lots of data on a lot of halos yet.
- We then use simulations which included ‘all/some’ physics (with lots of halos and halo properties).

Machine Learning

Populating Halos with Galaxies

- We want to use Machine Learning to learn about how many galaxies (or specific kinds of galaxies) are in halos with certain properties.
- We can do this with real data, however, we are not in the era where we have lots of data on a lot of halos *yet*.
- We then use simulations which included ‘all/some’ physics (with lots of halos and halo properties).

Machine Learning

Populating Halos with Galaxies

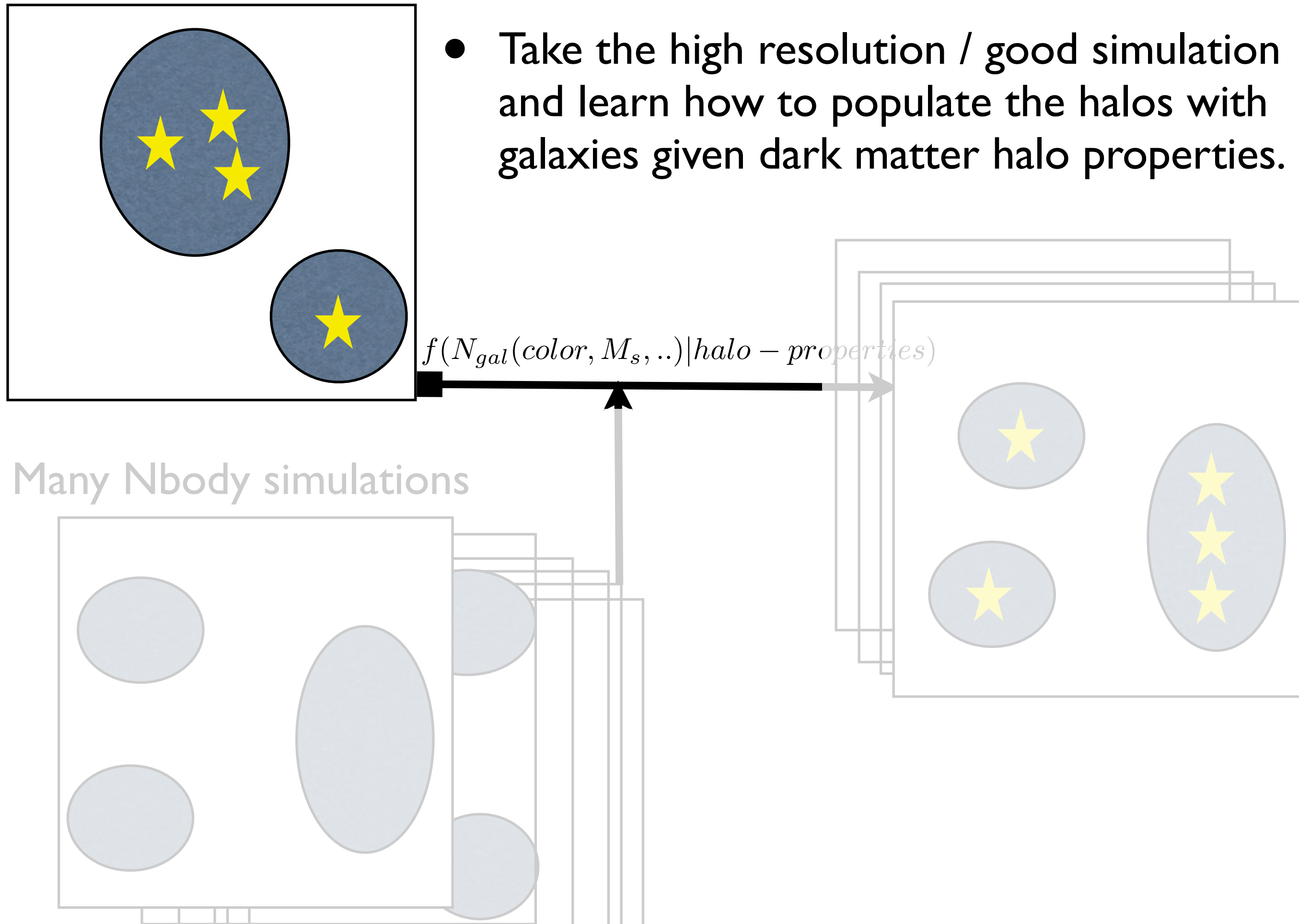
- We want to use Machine Learning to learn about how many galaxies (or specific kinds of galaxies) are in halos with certain properties.
- We can do this with real data, however, we are not in the era where we have lots of data on a lot of halos *yet*.
- We then use simulations which included ‘all/some’ physics (with lots of halos and halo properties).

Realistic simulations / Observations

- Take the high resolution / good simulation and learn how to populate the halos with galaxies given dark matter halo properties.

Many Nbody simulations

$$f(N_{gal}(color, M_s, ..) | halo - properties)$$

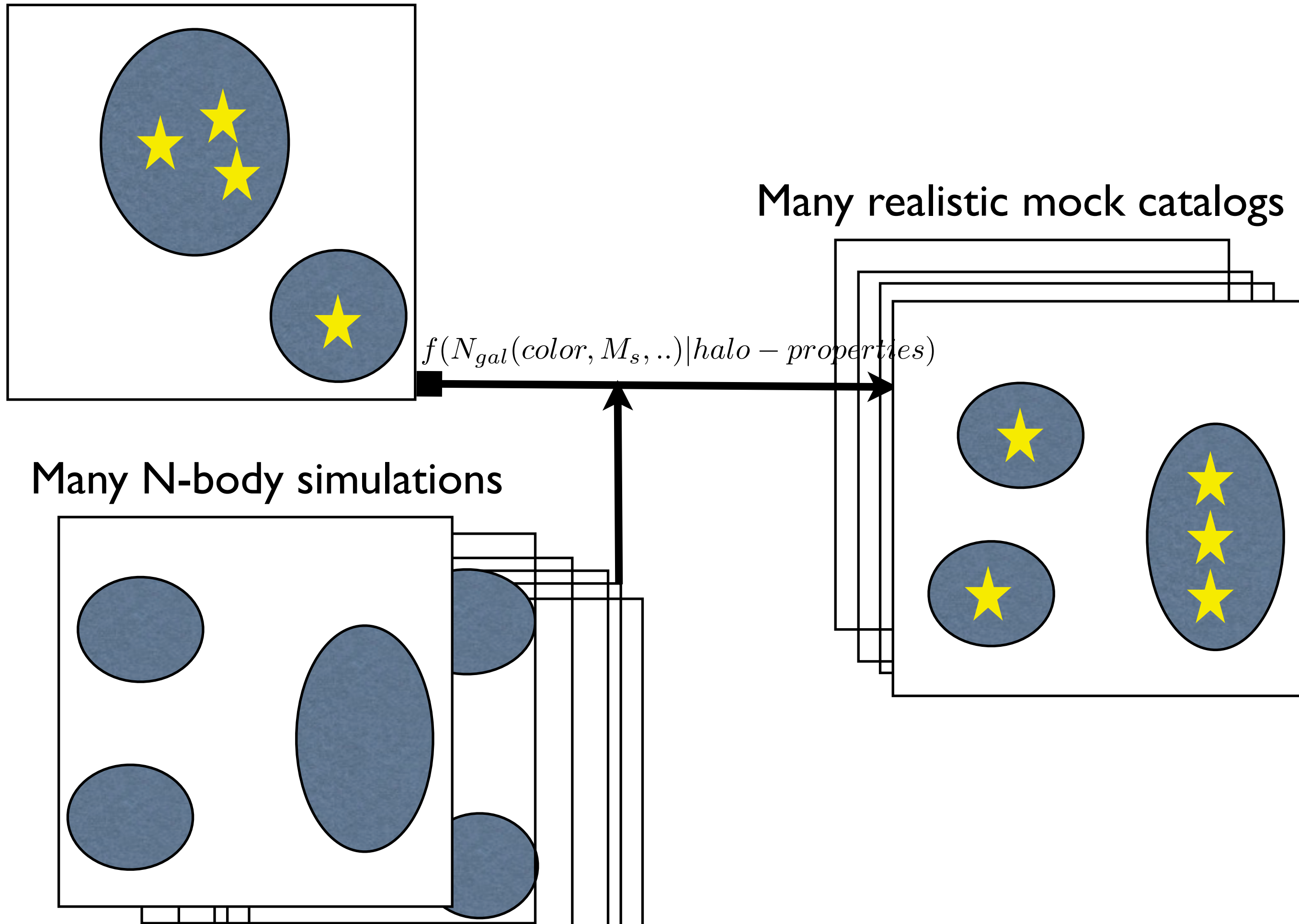


Few Realistic simulations / Observations

Many realistic mock catalogs

$$f(N_{gal}(color, M_s, ..) | halo - properties)$$

Many N-body simulations

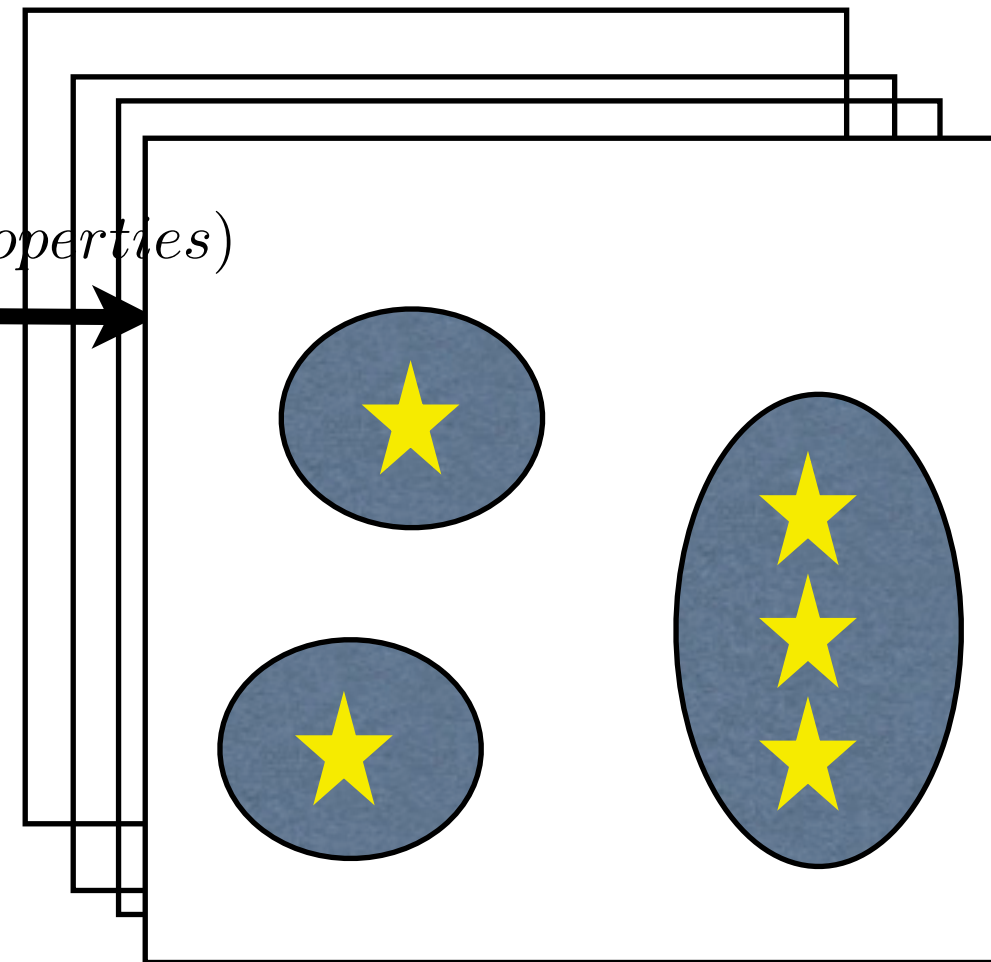
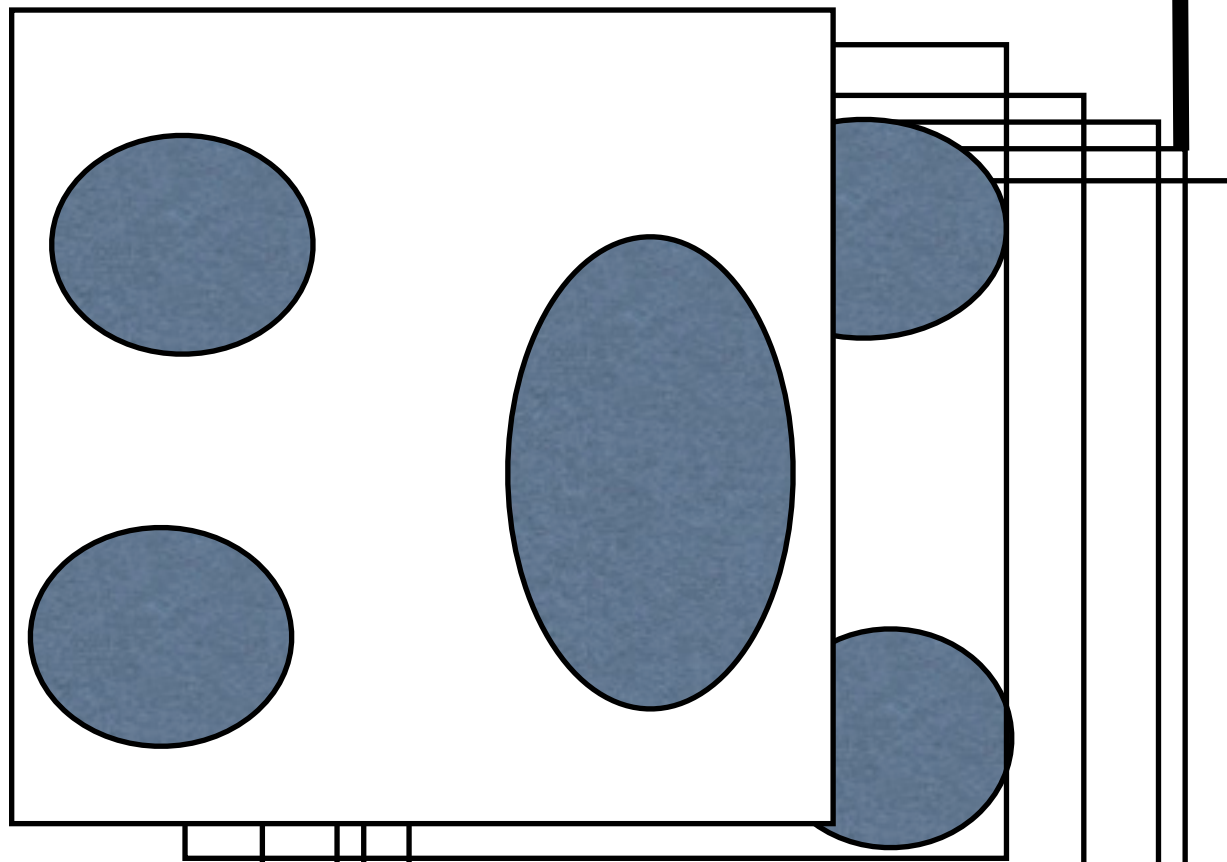


Realistic simulations / Observations

- Can then take dark matter halos from large number of not-as-expensive DM only simulations, and make many independent mock catalogs.

$$f(N_{gal}(color, M_s, ..)| halo - properties)$$

Many Nbody simulations



Not enough observations... so we use Millenium

- We use the halo catalogues and semi-analytic galaxies from the Millennium simulation: $\Omega_m=0.25$, $\Omega_b=0.045$, $\Omega_\Lambda=0.75$, $h=0.73$, $n_s=1.0$, $\sigma_8=0.9$.
- We only use the central and satellite galaxies of the primary Millennium halos with mass $> 10^{12}M_{\text{sun}}/h$.
- There are about 400,000 of these halos. We subsample to 60,000 for some tests.
- We split the sample/subsample randomly and equally into training and test sets.

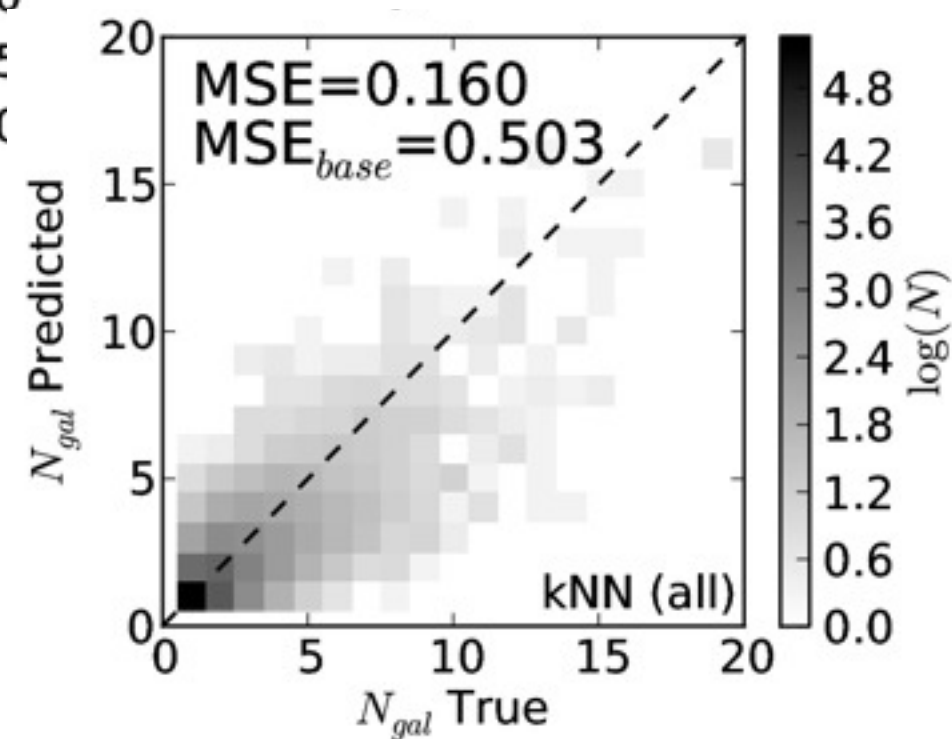
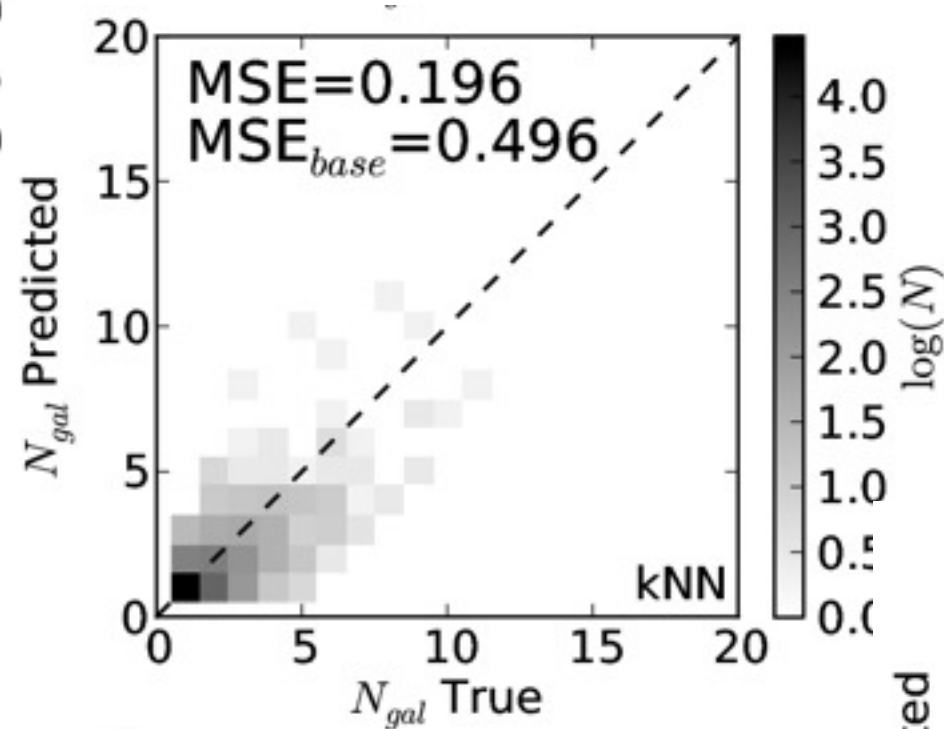
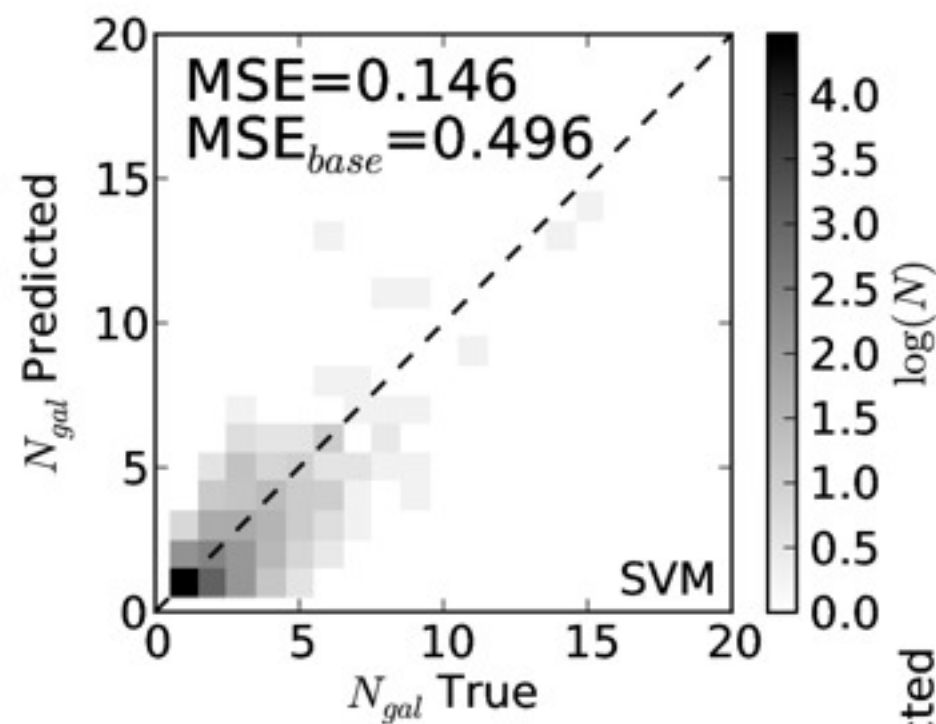
Process

- We split the data into 3 sets
- One for training
- One for validation
- One for testing (we predict the number of galaxies given what we learn from the training set)

Results - Number of galaxies

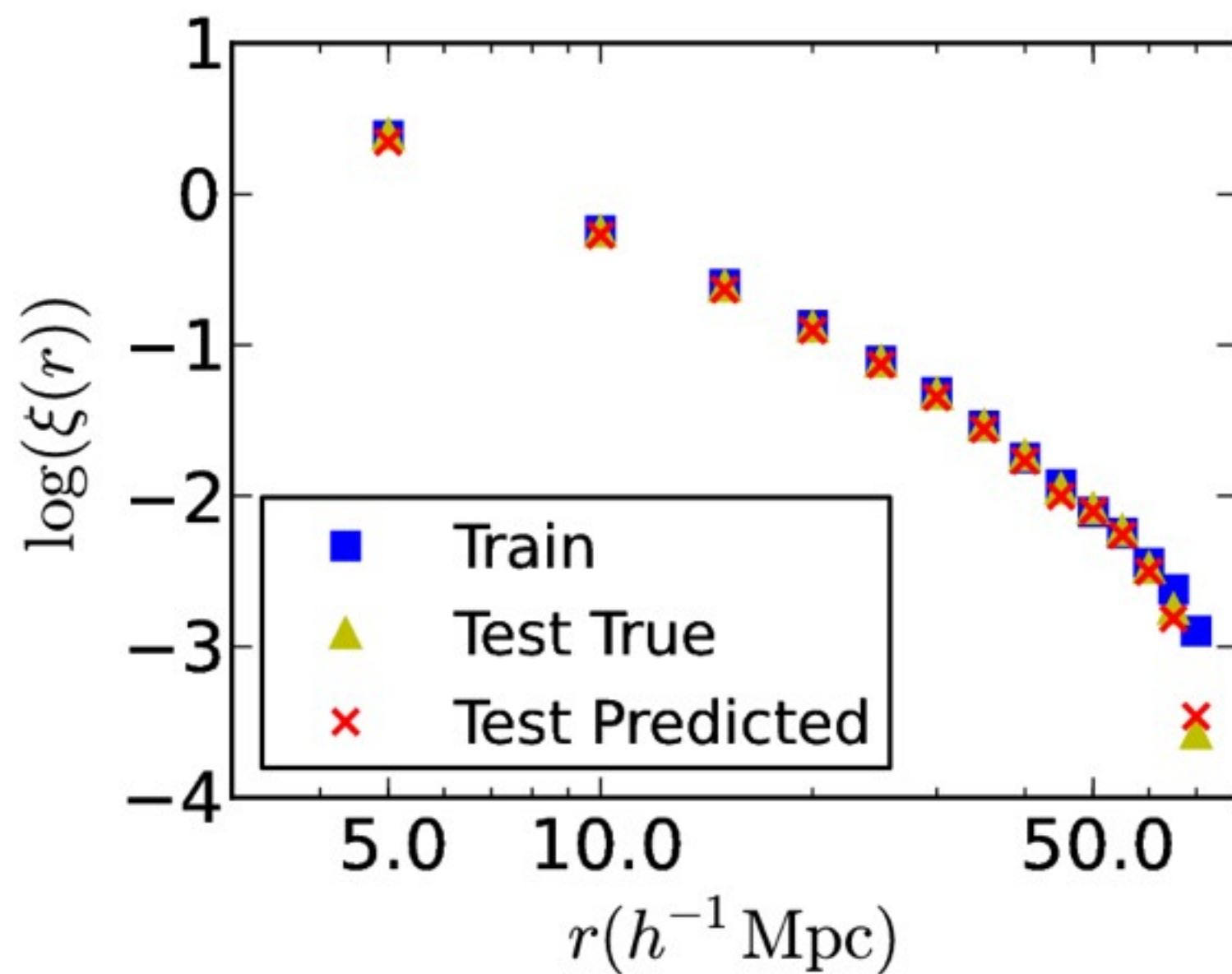
Xu, S.H, Trac, Schneider, Poczos 2013

Results - Number of galaxies



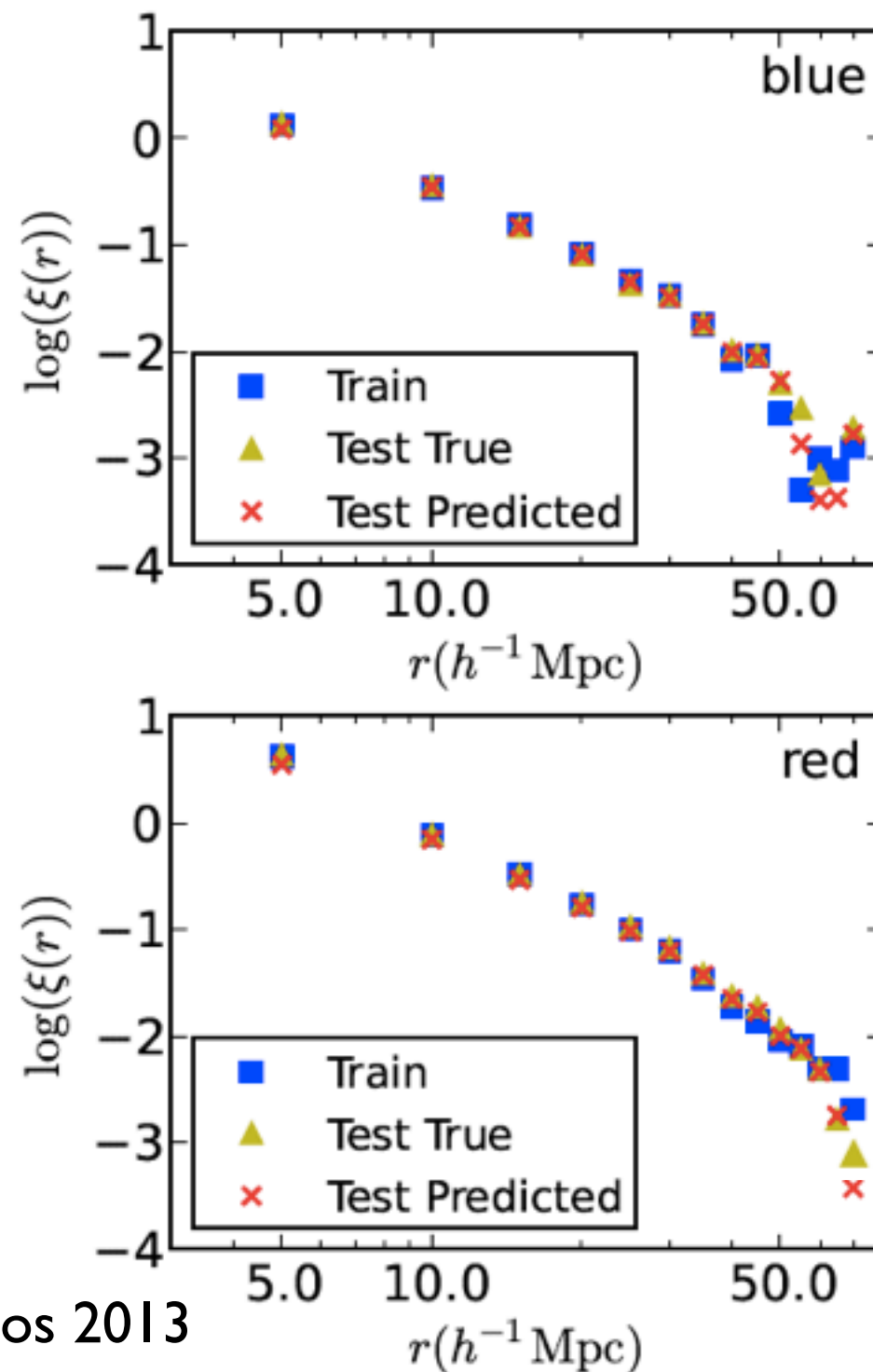
Xu, S.H, Trac, Schneider, Poczos 2013

Results - Correlation functions



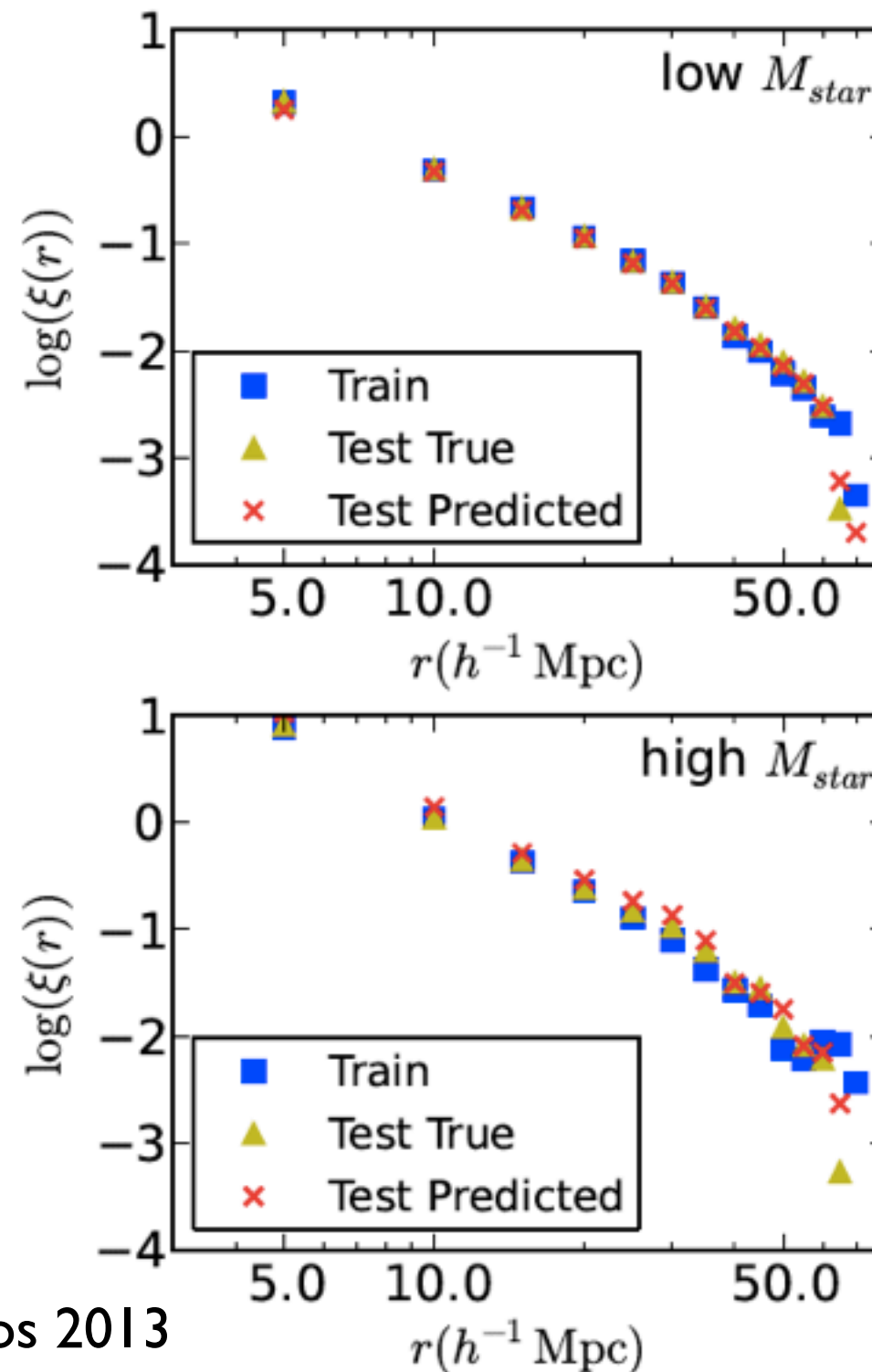
Xu, S.H, Trac, Schneider, Poczós 2013

Results - Correlation functions of blue and red galaxies

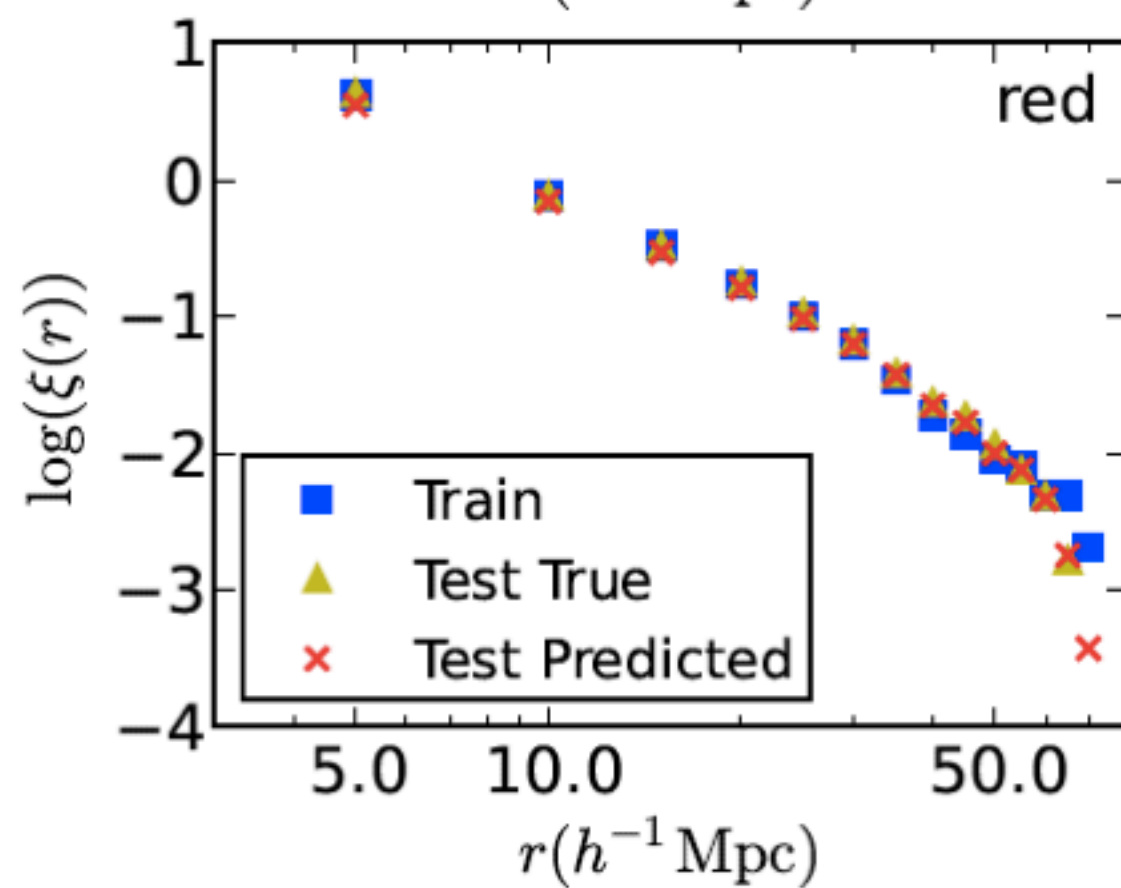
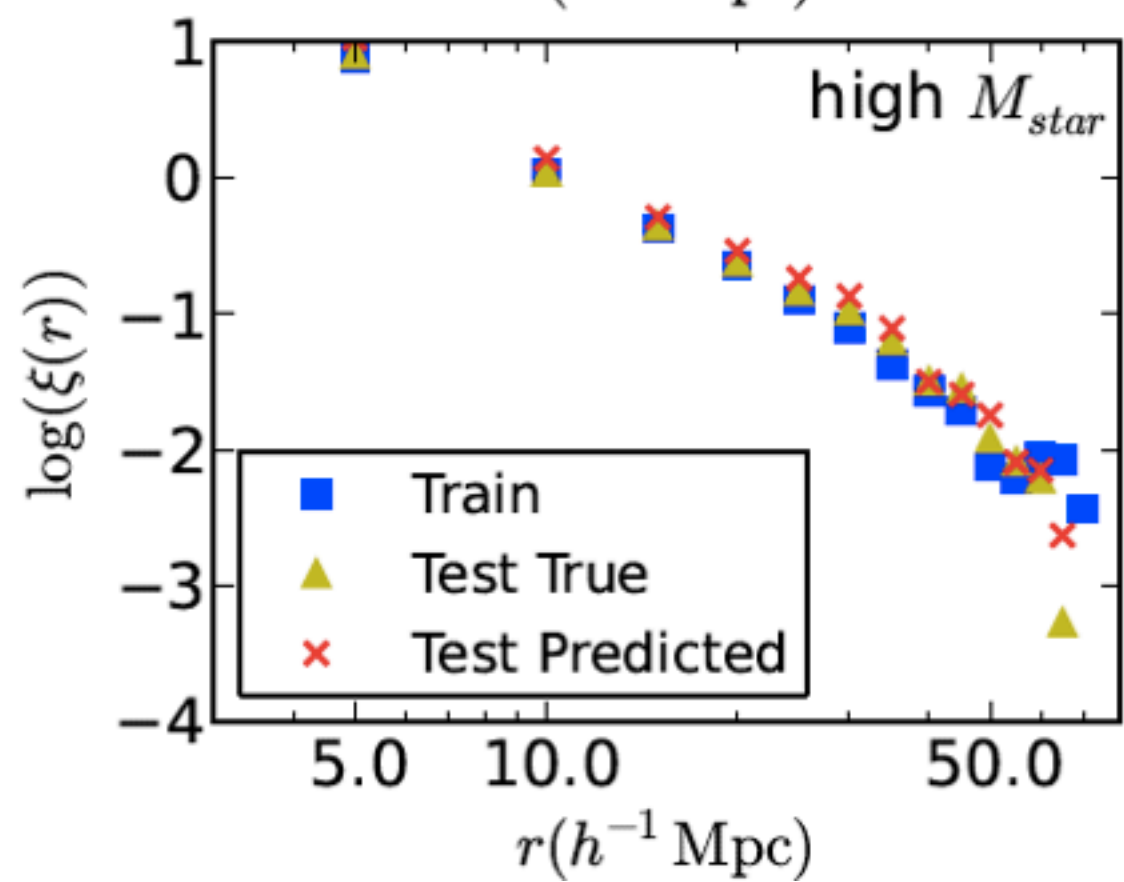
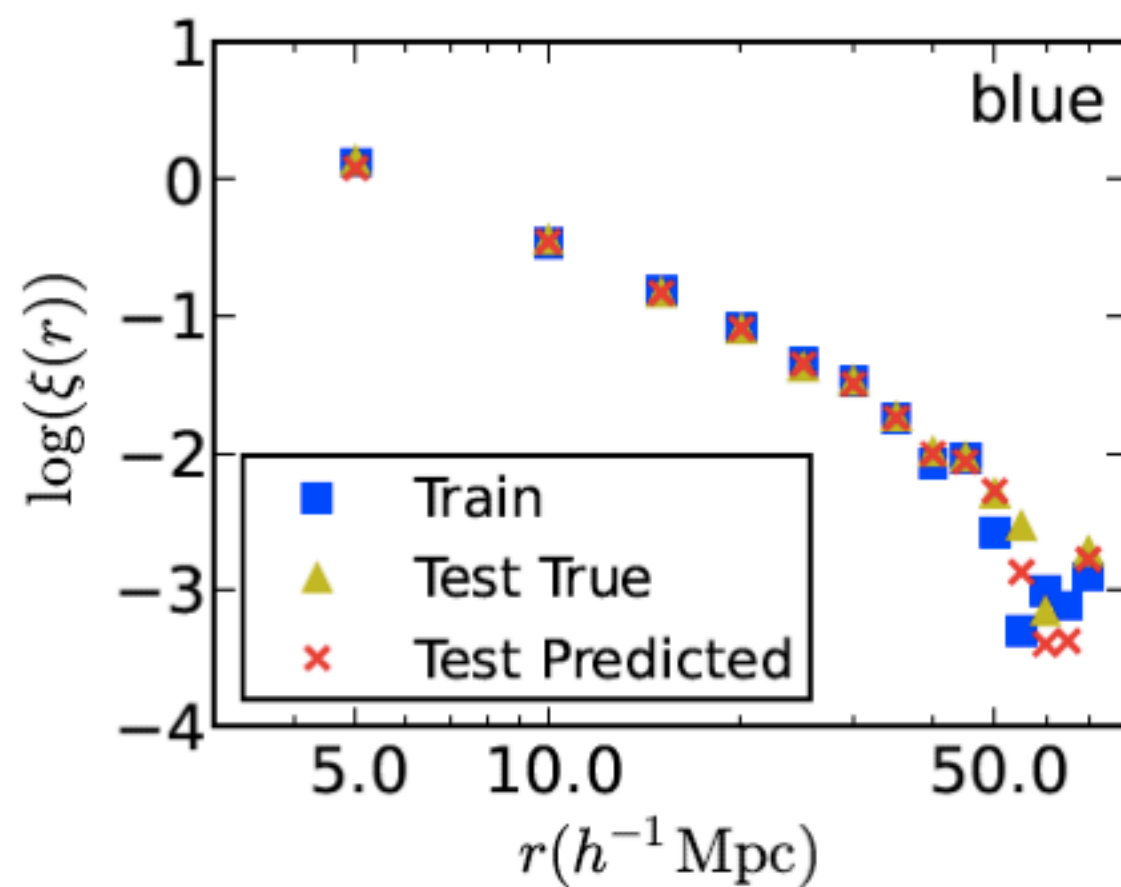
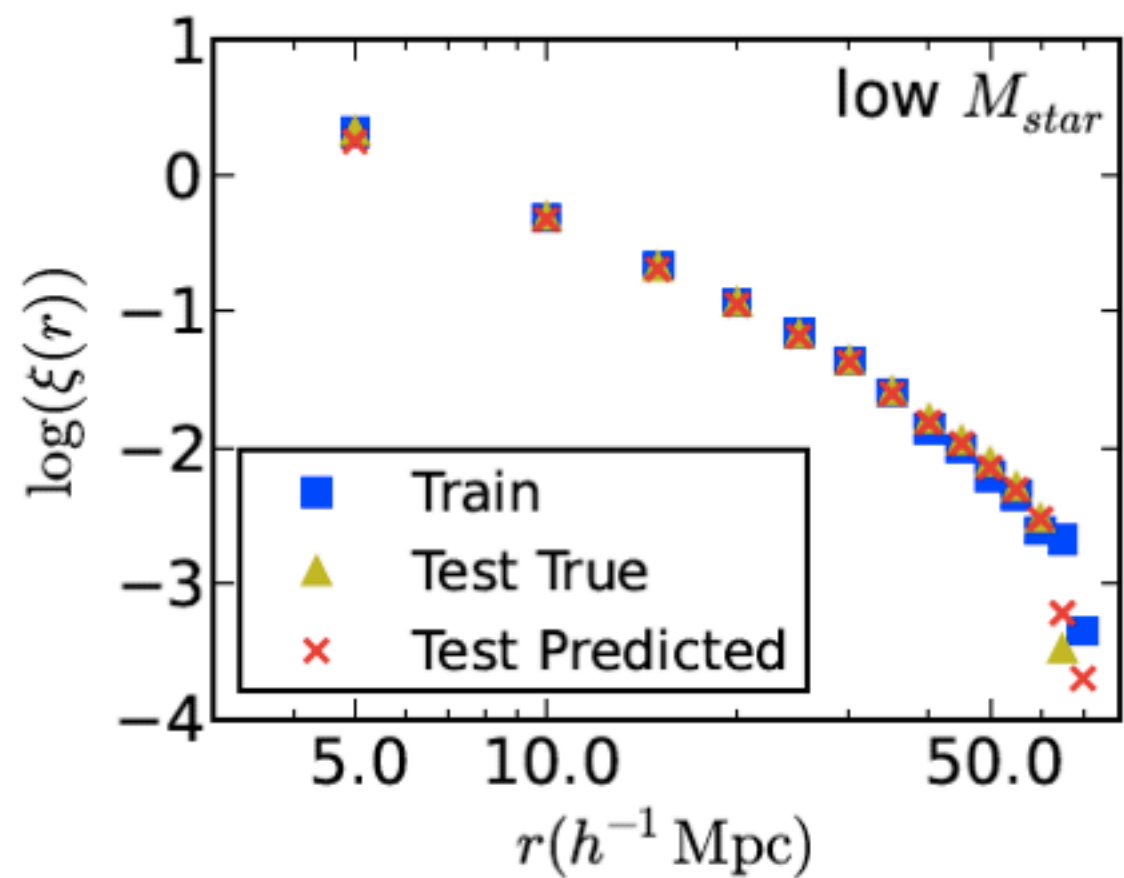


Xu, S.H, Trac, Schneider, Poczos 2013

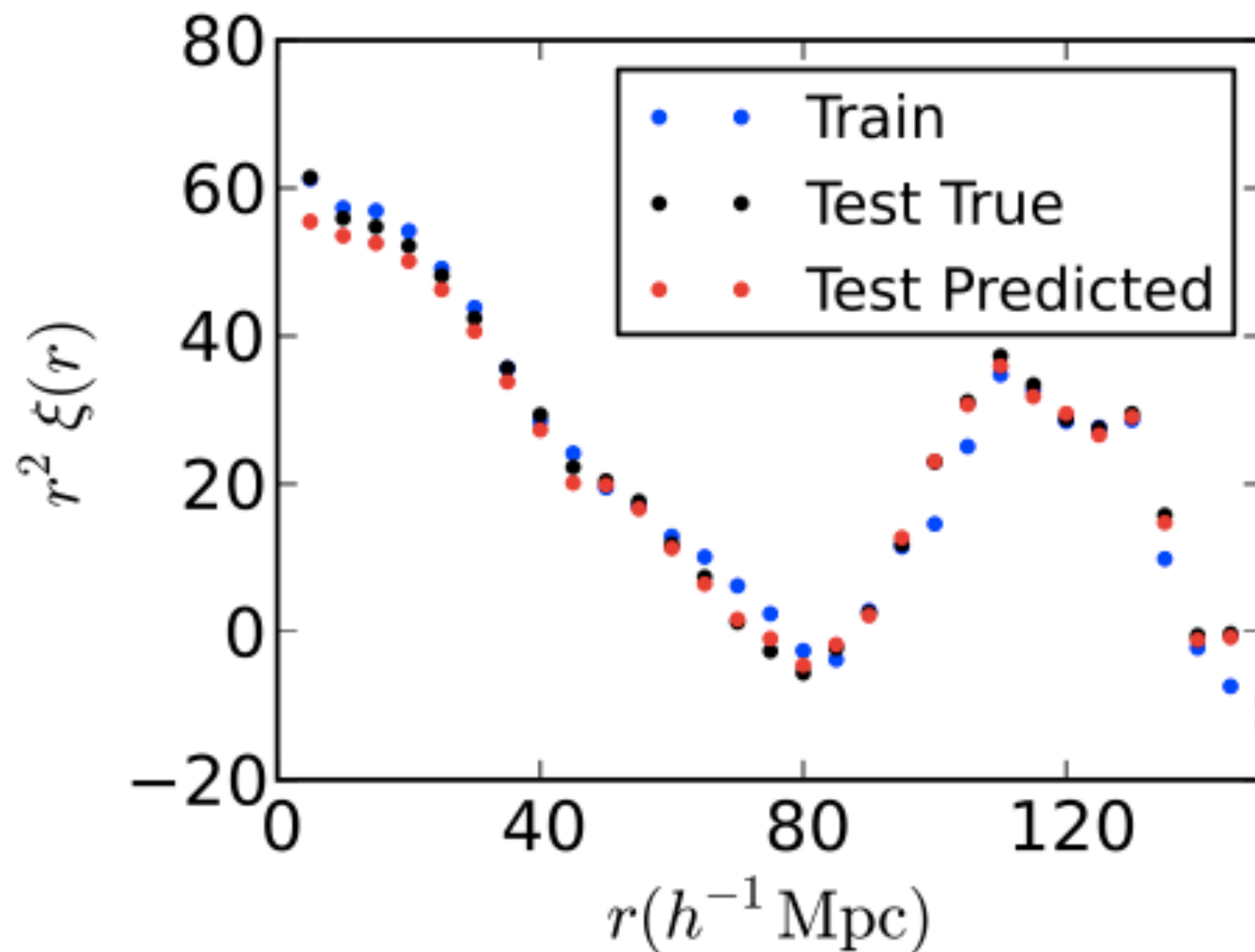
Results - Correlation functions of galaxies with different stellar mass thresholds



Xu, S.H, Trac, Schneider, Poczos 2013



Results VI - Correlation function at large scale (BAO)



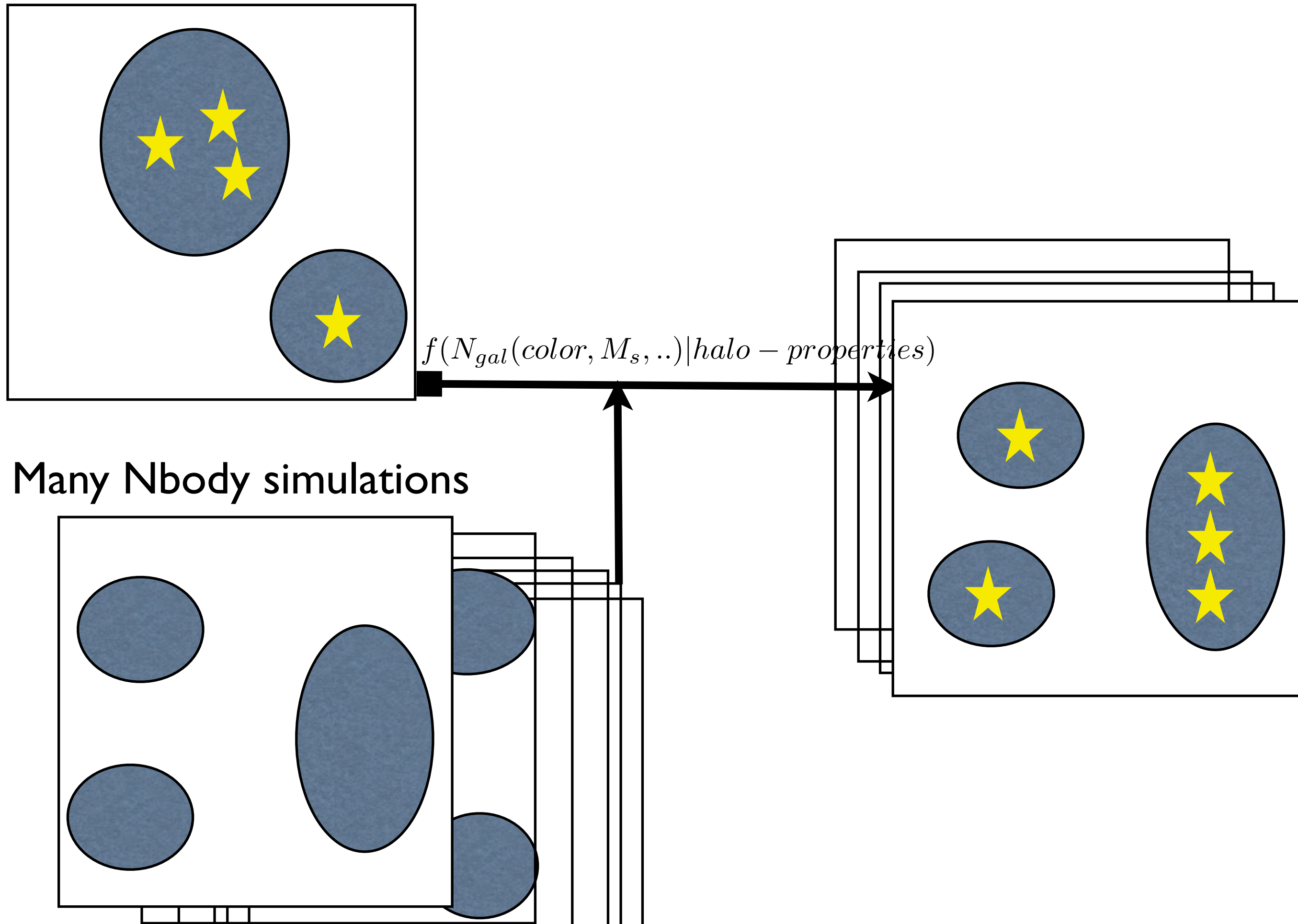
Note this is using only 60,000 galaxies,
so there is quite a bit of shot noise.

Xu, S.H, Trac, Schneider, Poczos 2013

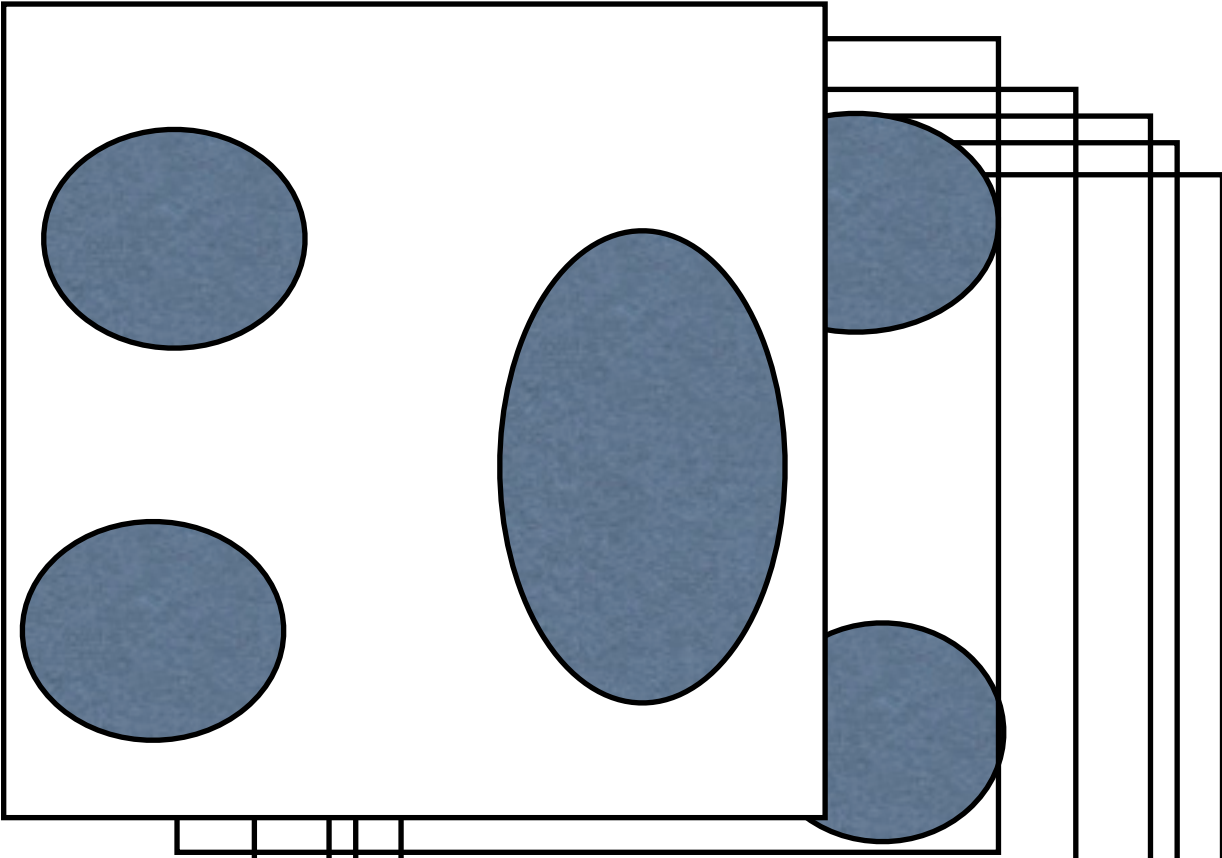
Mini-Conclusion

- Machine learning offers a completely **model-independent** method for understanding the halo-galaxy mapping while avoiding subhalo finding.
- ML techniques give robust predictions of the number of galaxies per halo, the distribution of halos with N_{gal} and the galaxy correlation function.
- **Now, we move to even more ambitious question:**
 - **Can we skip the whole process of N-body simulations with Machine Learning?**

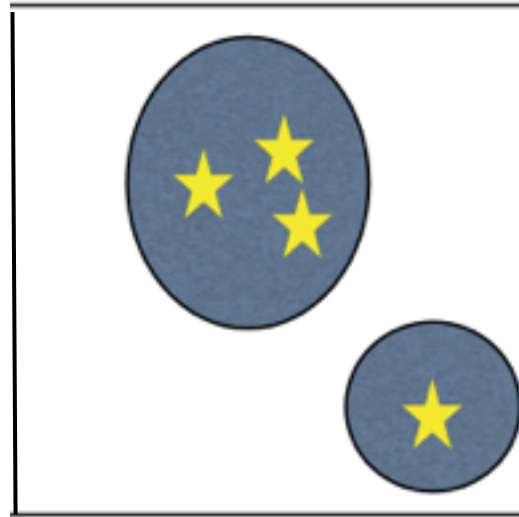
Realistic simulations / Observations



Many Nbody simulations

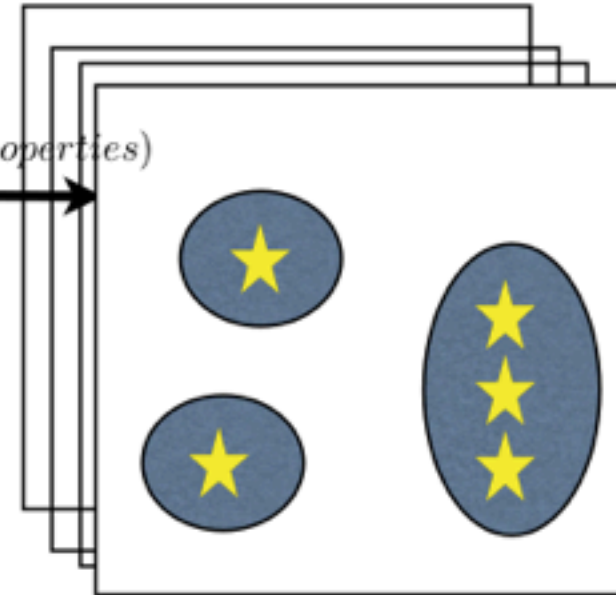
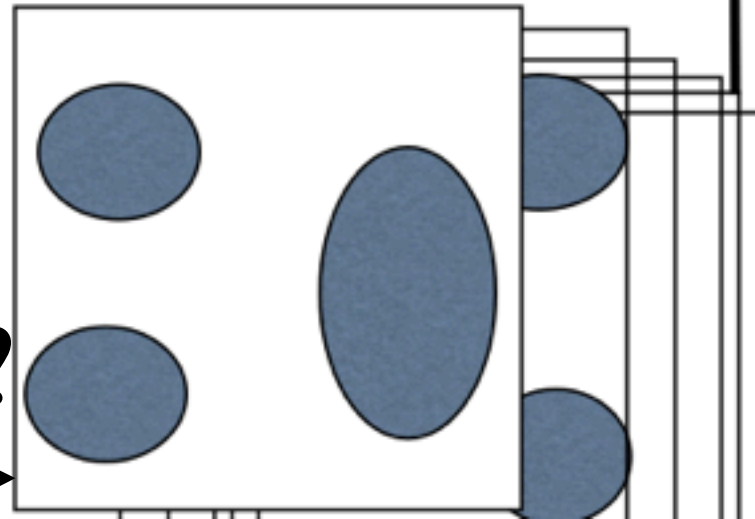


Realistic simulations / Observations

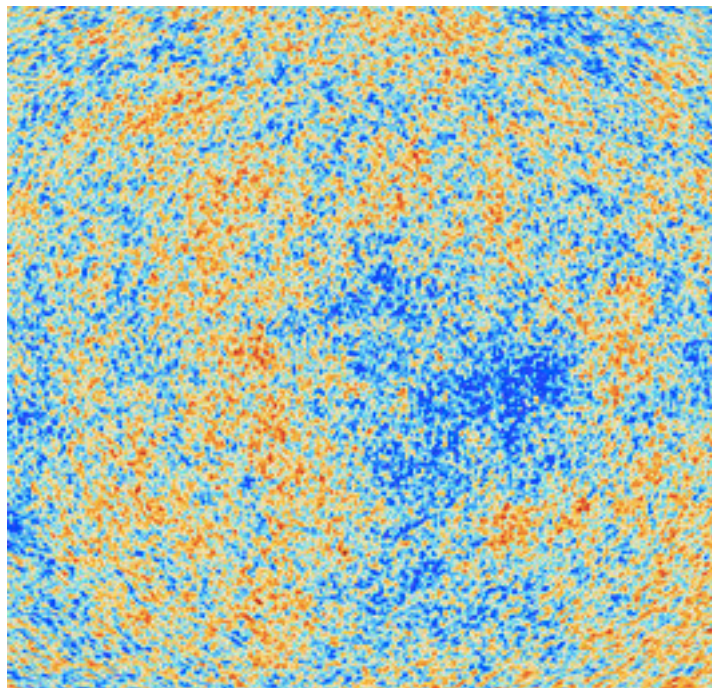


$f(N_{gal}(color, M_s, ..)| halo - properties)$

Many Nbody simulations



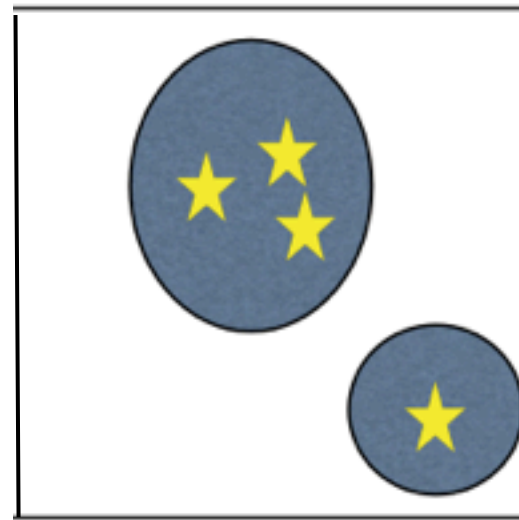
Machine learning?



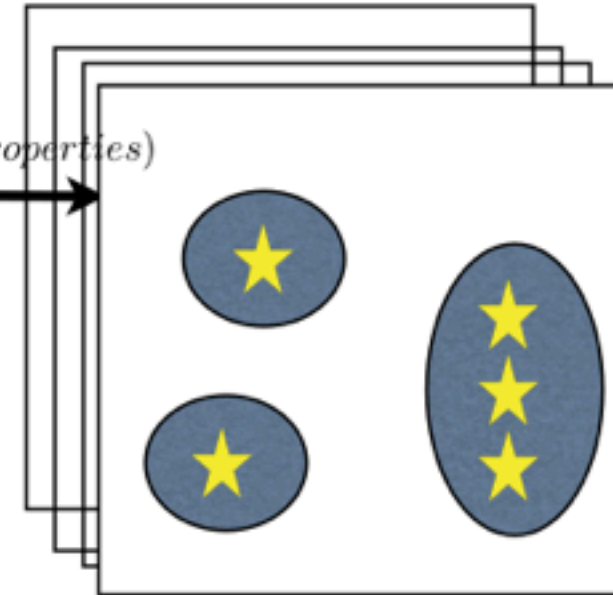
at $z=0.5$

at $z=1100$

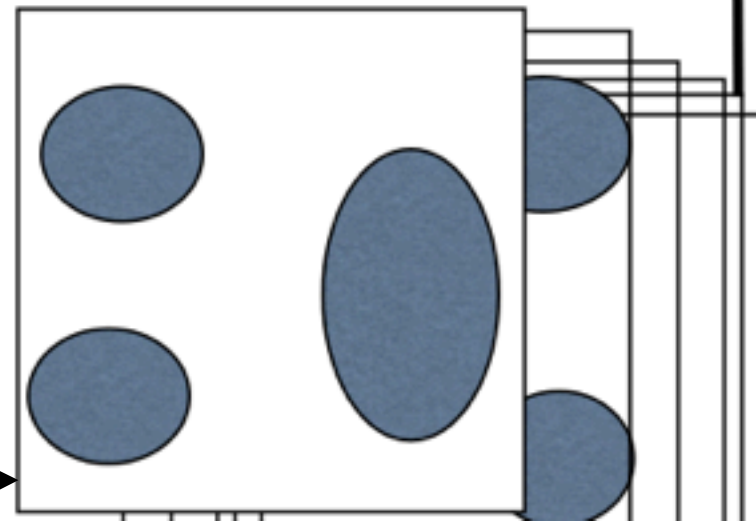
Realistic simulations / Observations



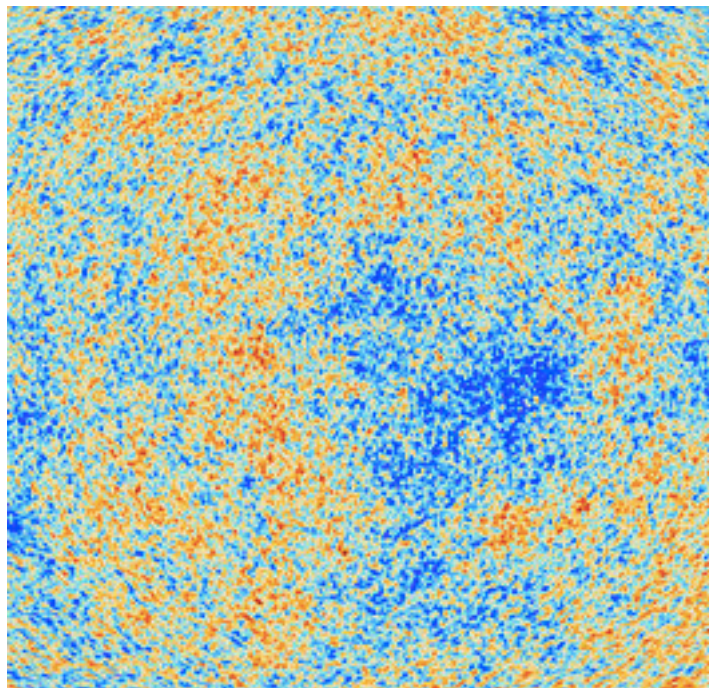
$f(N_{gal}(color, M_s, ..)| halo - properties)$



Many Nbody simulations



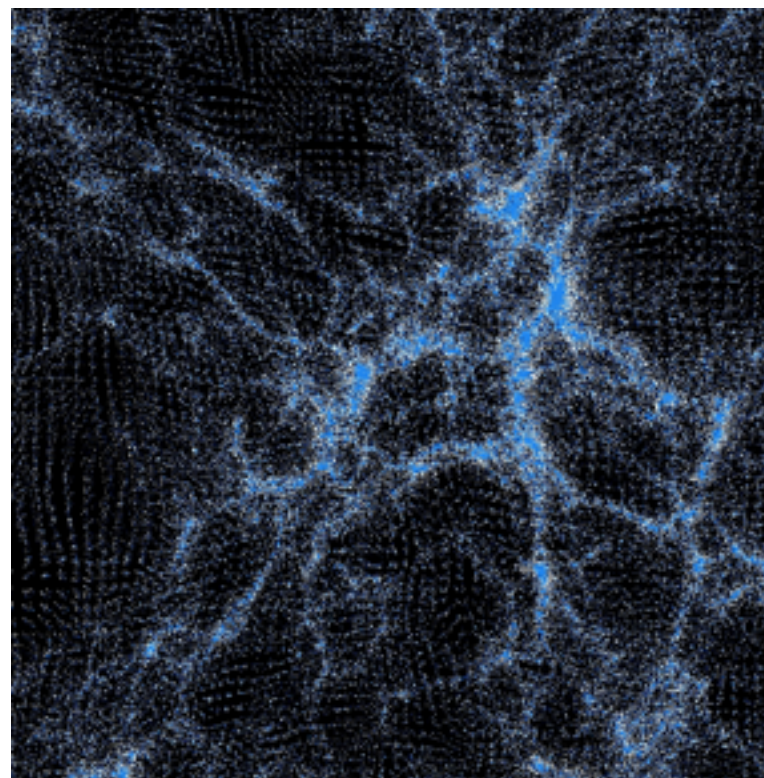
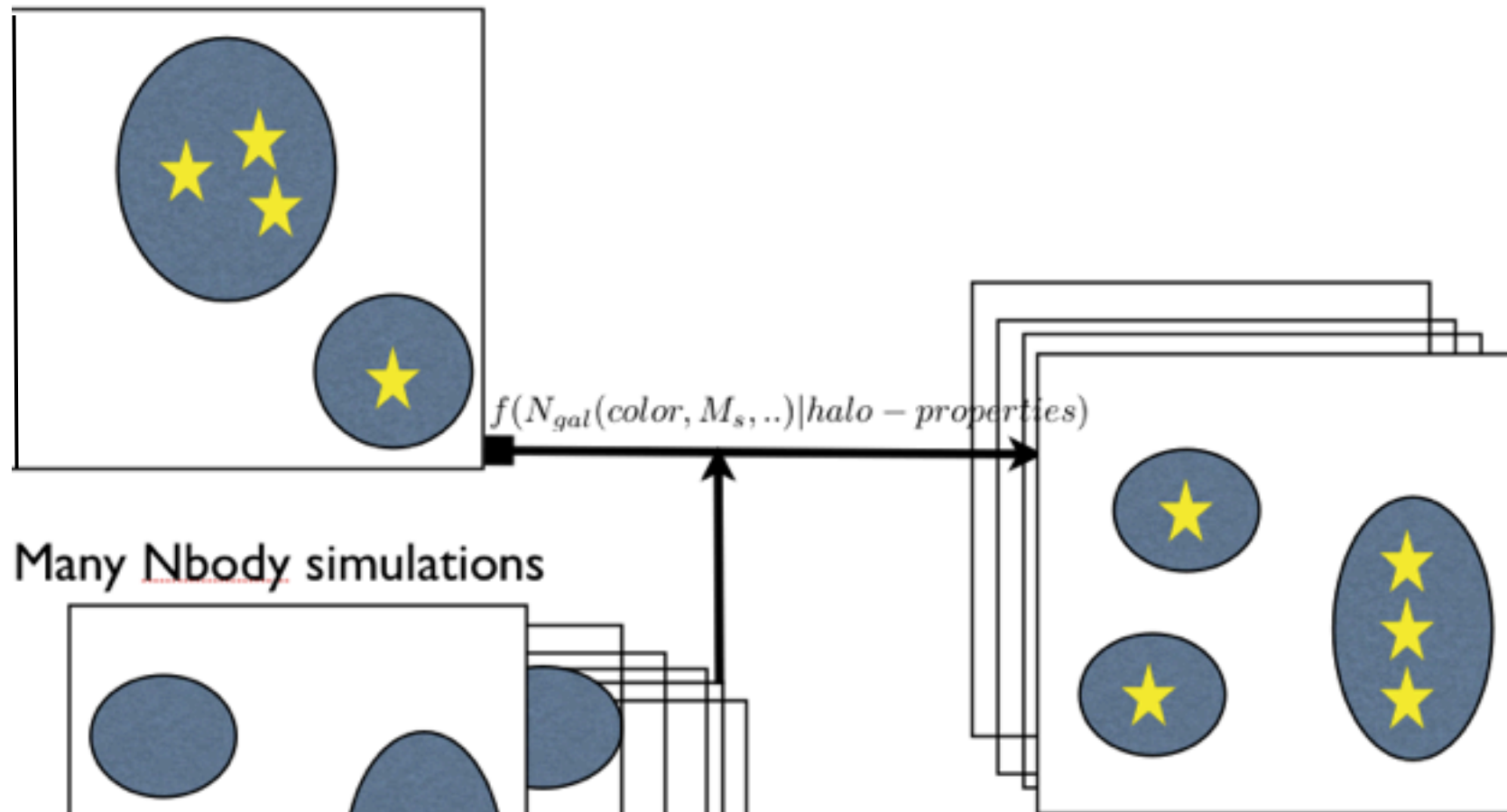
at $z=0.5$



at $z=1100$

New method in machine learning allows
distribution-to-distribution prediction

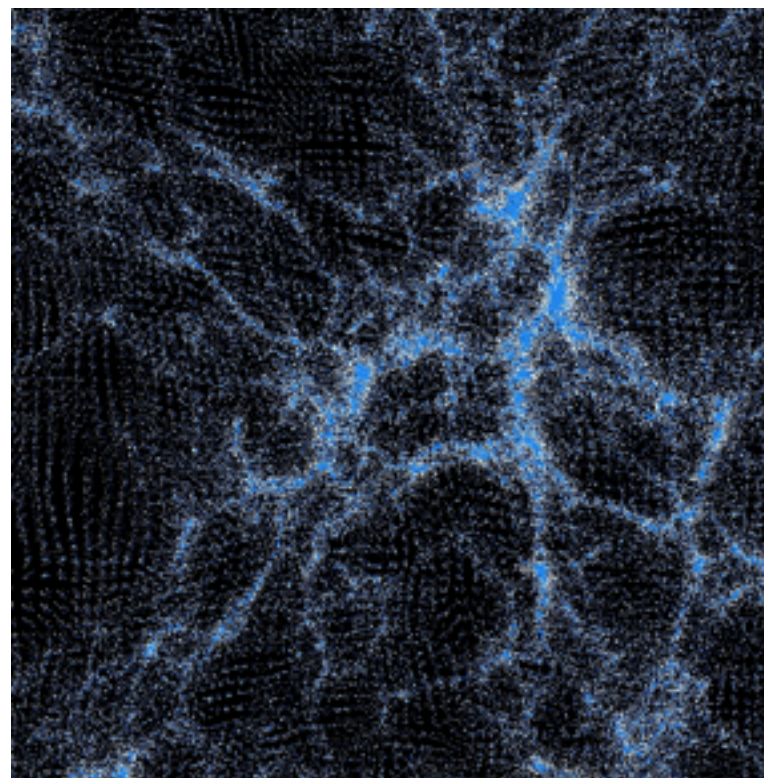
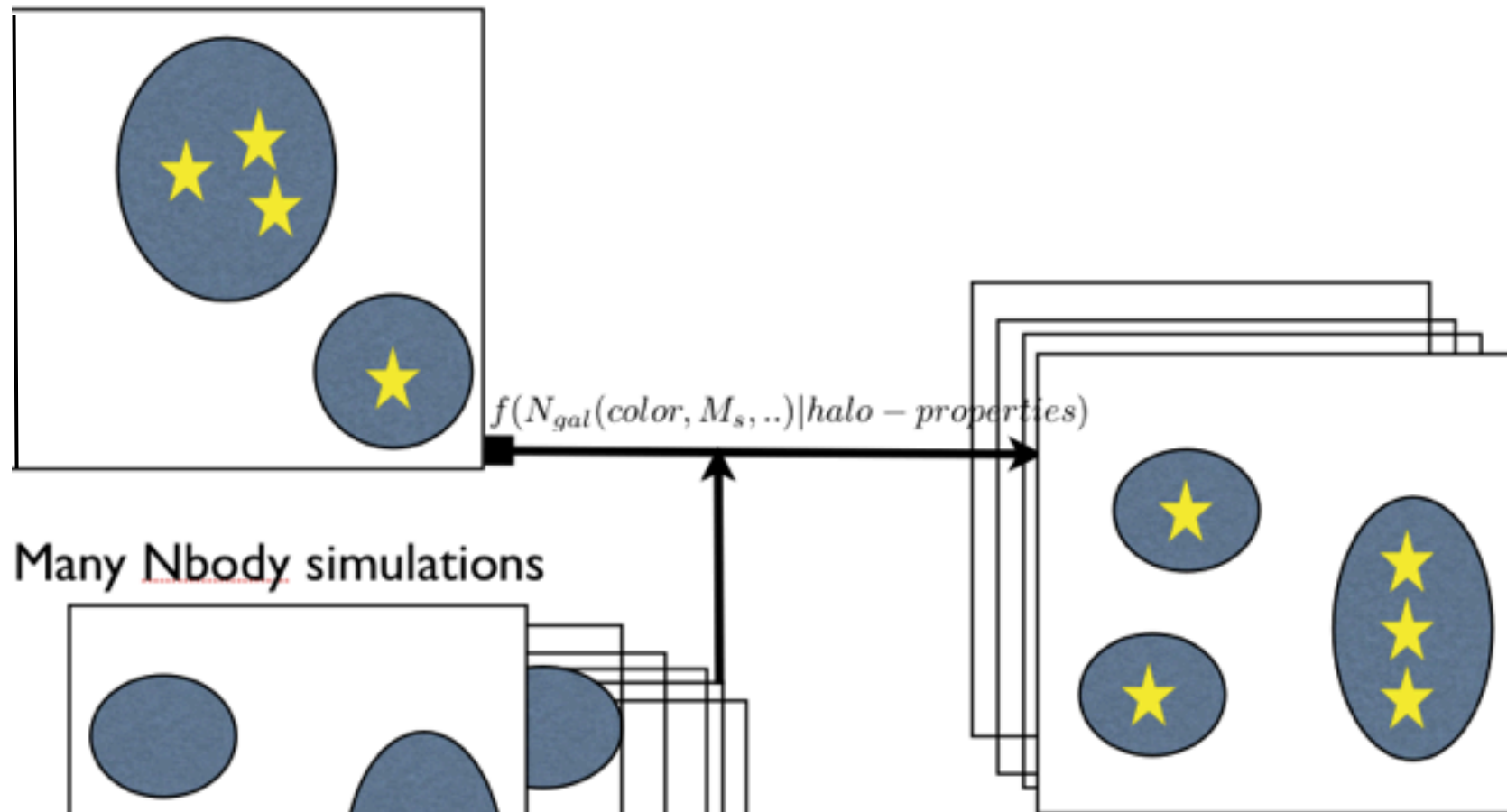
Realistic simulations / Observations



To make this easier, we evolve the field from initial conditions to $z=0.5$ using 2LPT (\sim mins)

2LPT to $z=0.5$

Realistic simulations / Observations



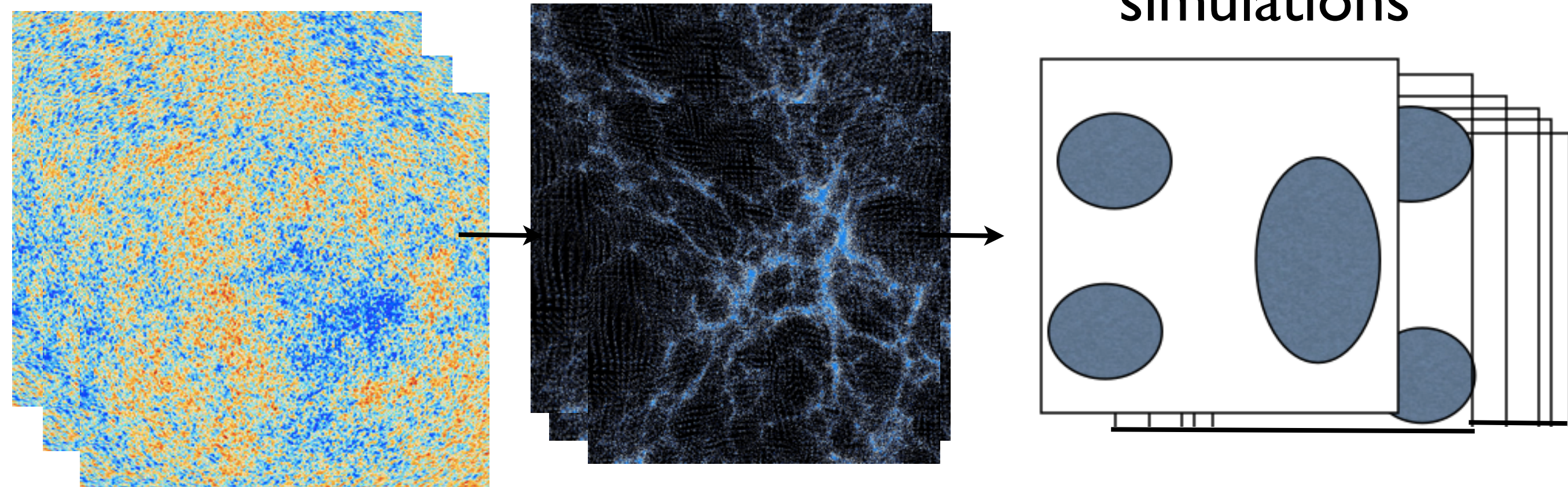
Now the question is:
can we learn how to evolve from the
2LPT field to Nbody density field?

2LPT to $z=0.5$

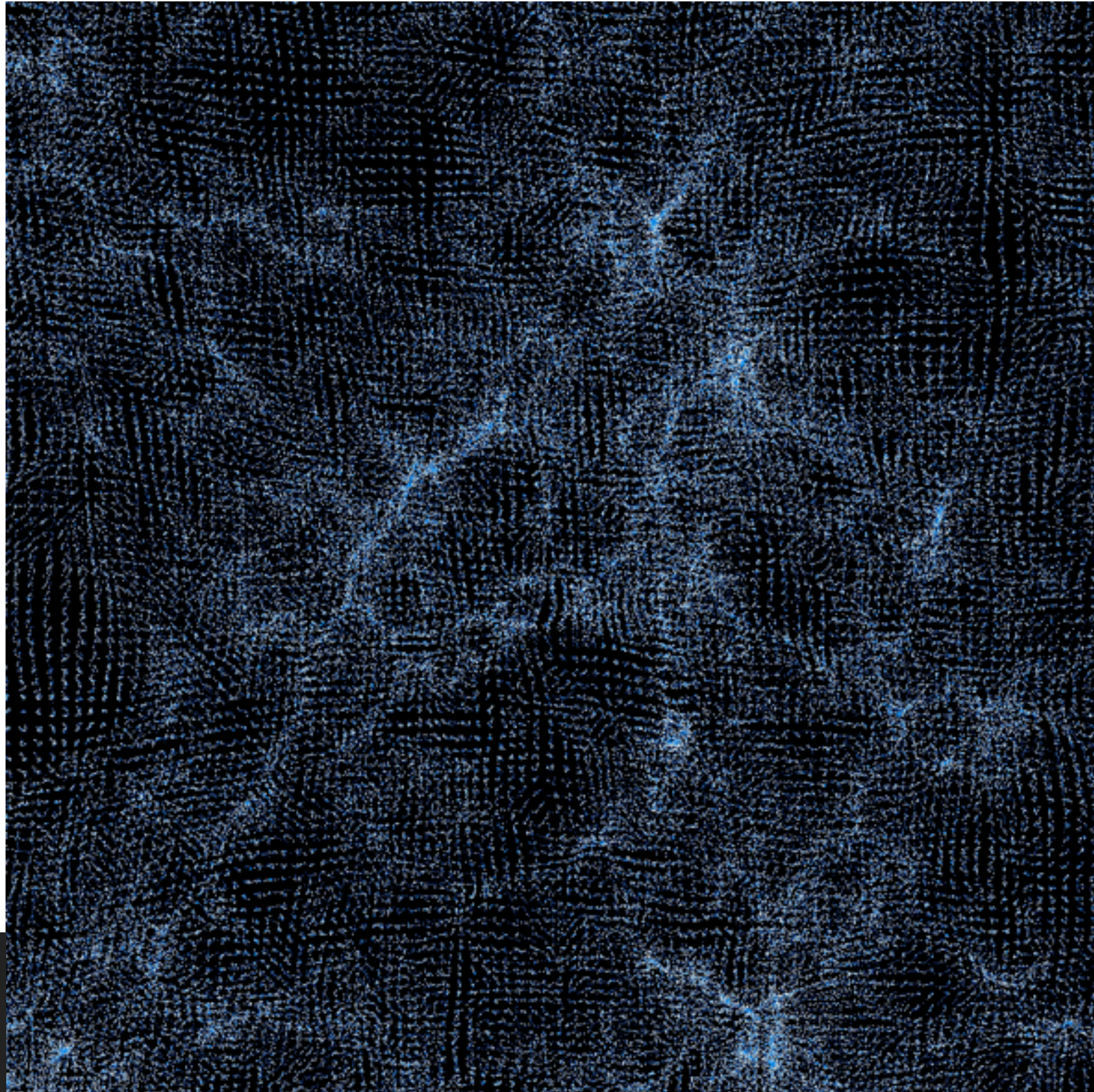
multiple initial
conditions

2LPT to $z=0.5$

Multiple near
N-body
simulations

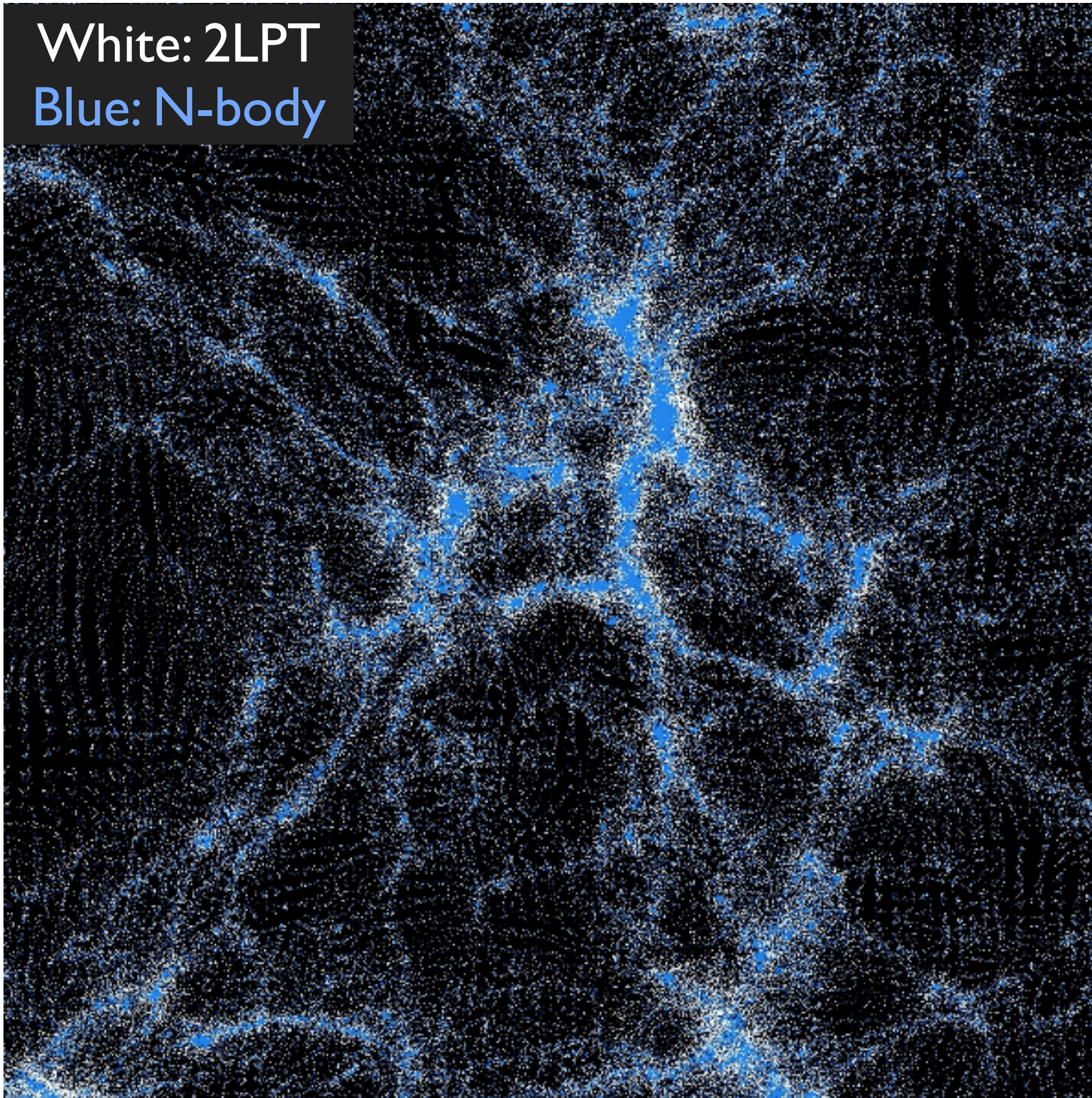


How does 2LPT field compares to N-body field?

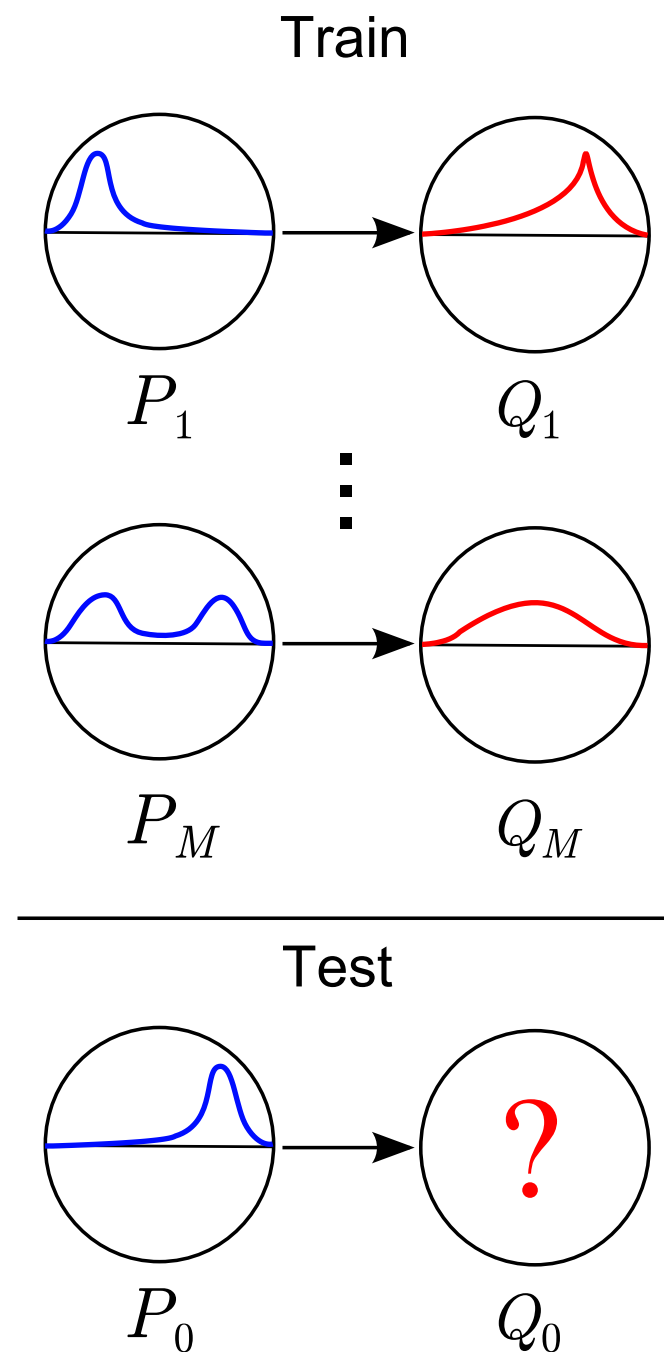


White: 2LPT
Blue: N-body

White: 2LPT
Blue: N-body

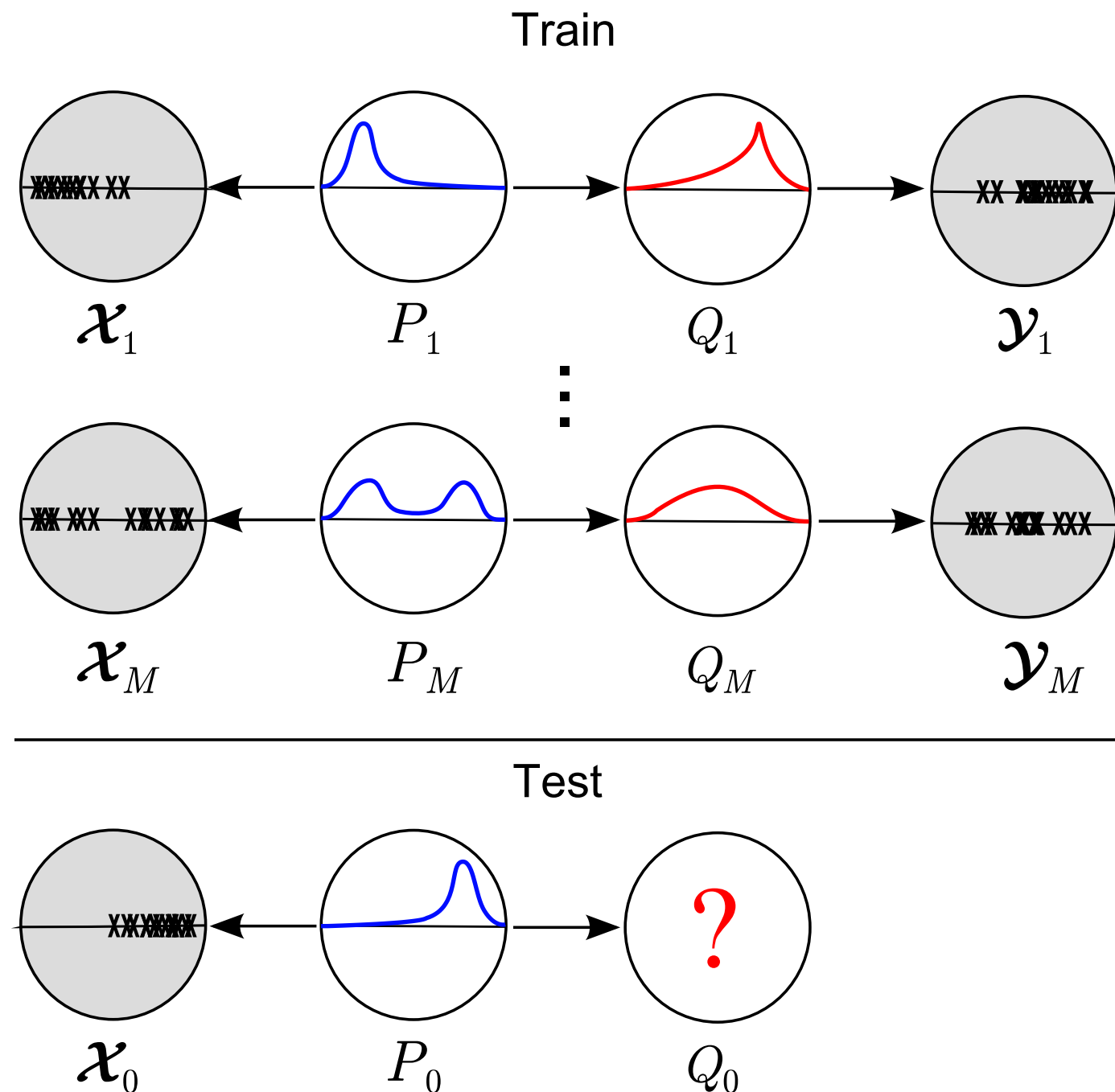


(2LPT) Distribution to (N-body) Distribution Machine Learning algorithms



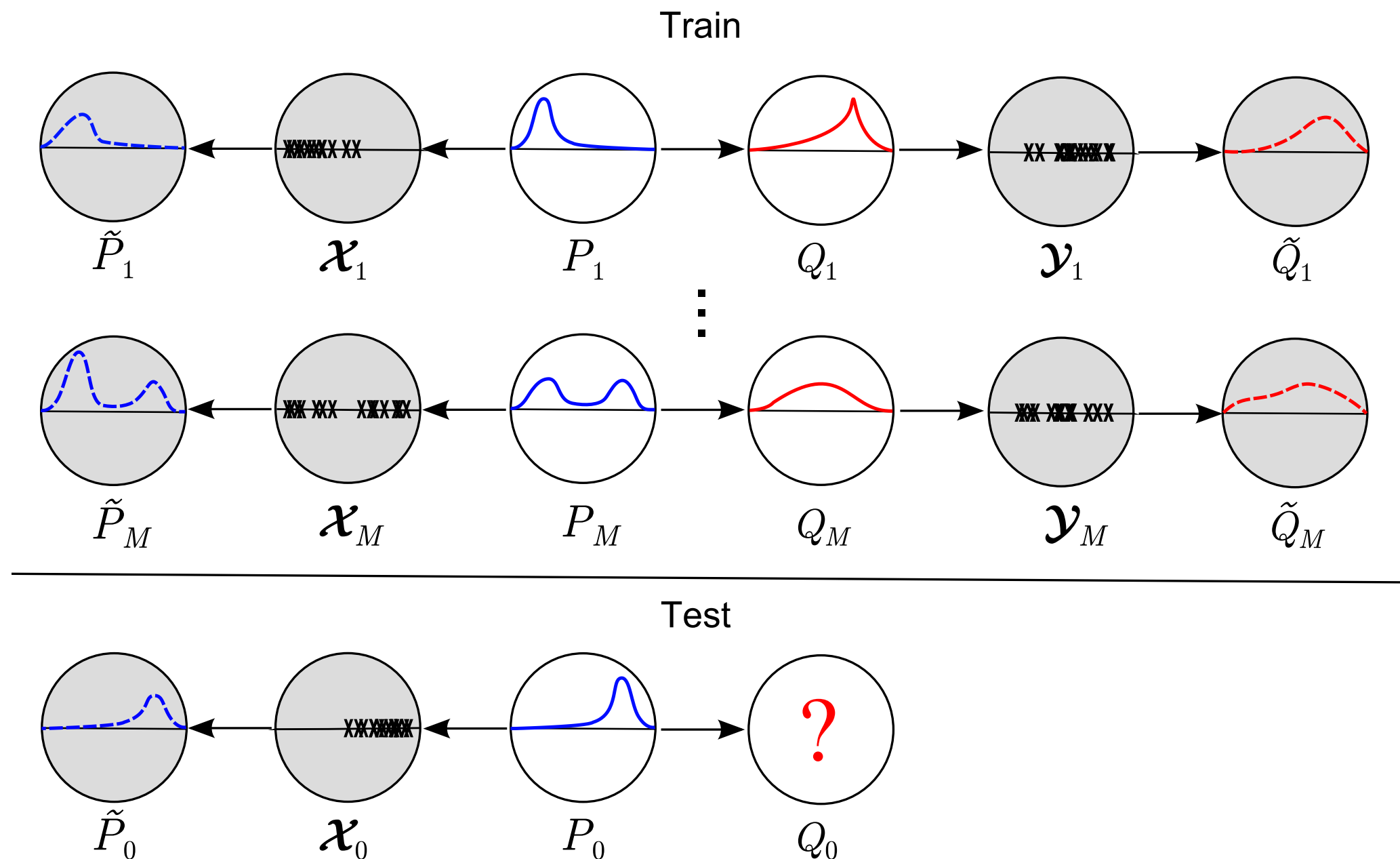
courtesy Junier Oliver

(2LPT) Distribution to (N-body) Distribution Machine Learning algorithms



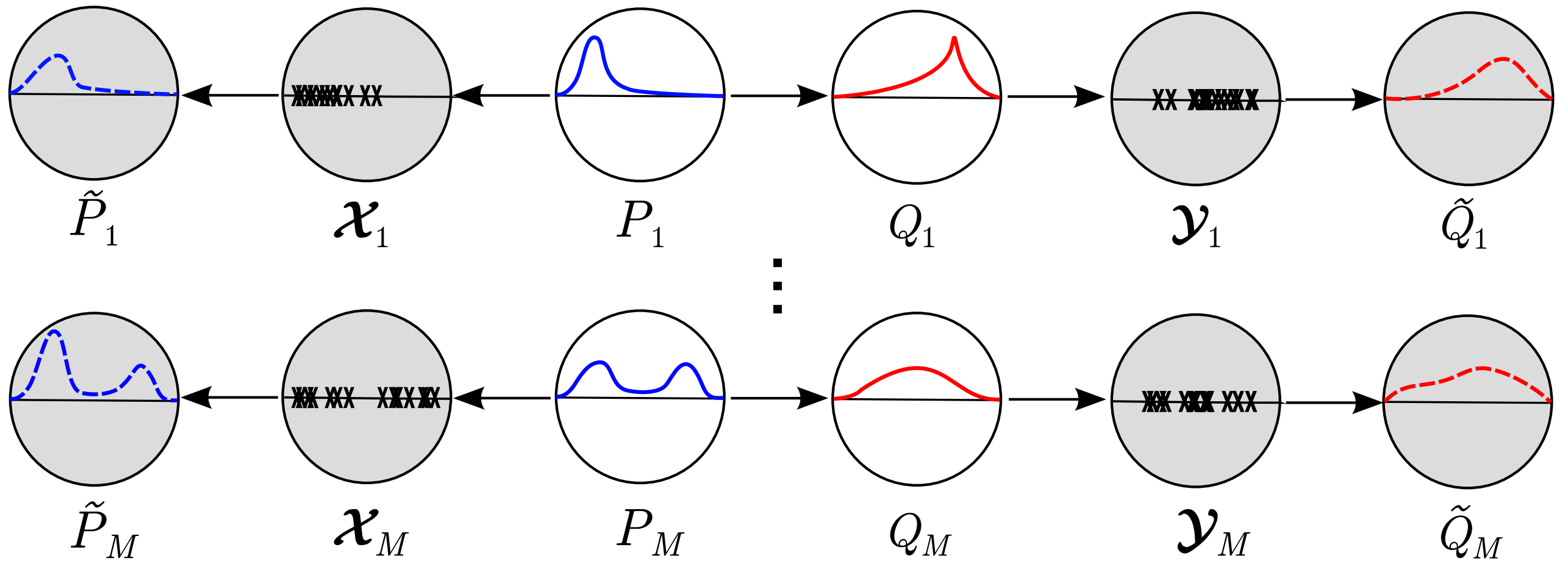
courtesy Junier Oliver

(2LPT) Distribution to (N-body) Distribution Machine Learning algorithms

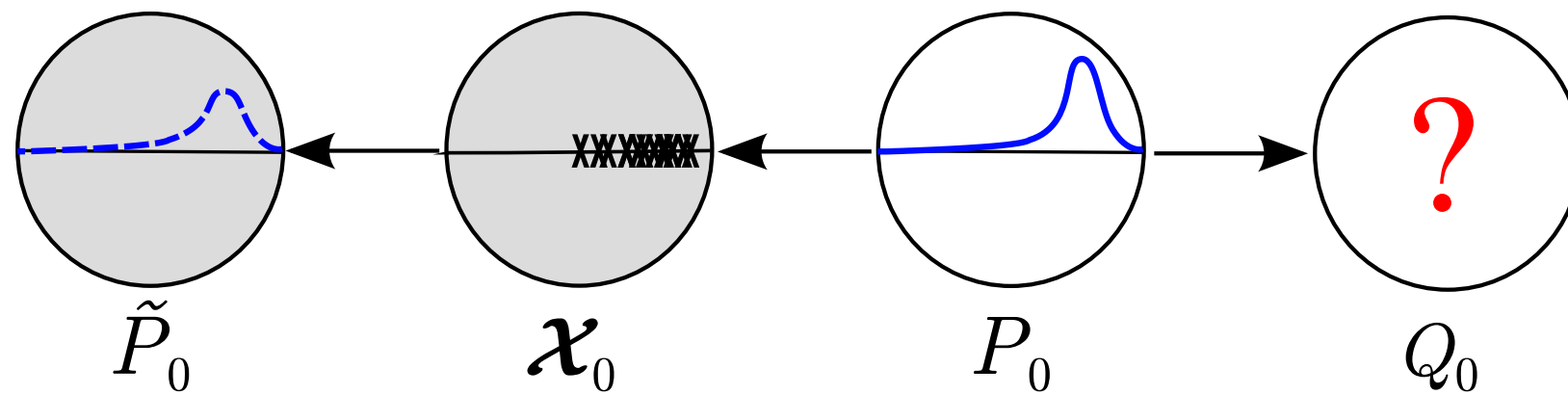


courtesy Junier Oliver

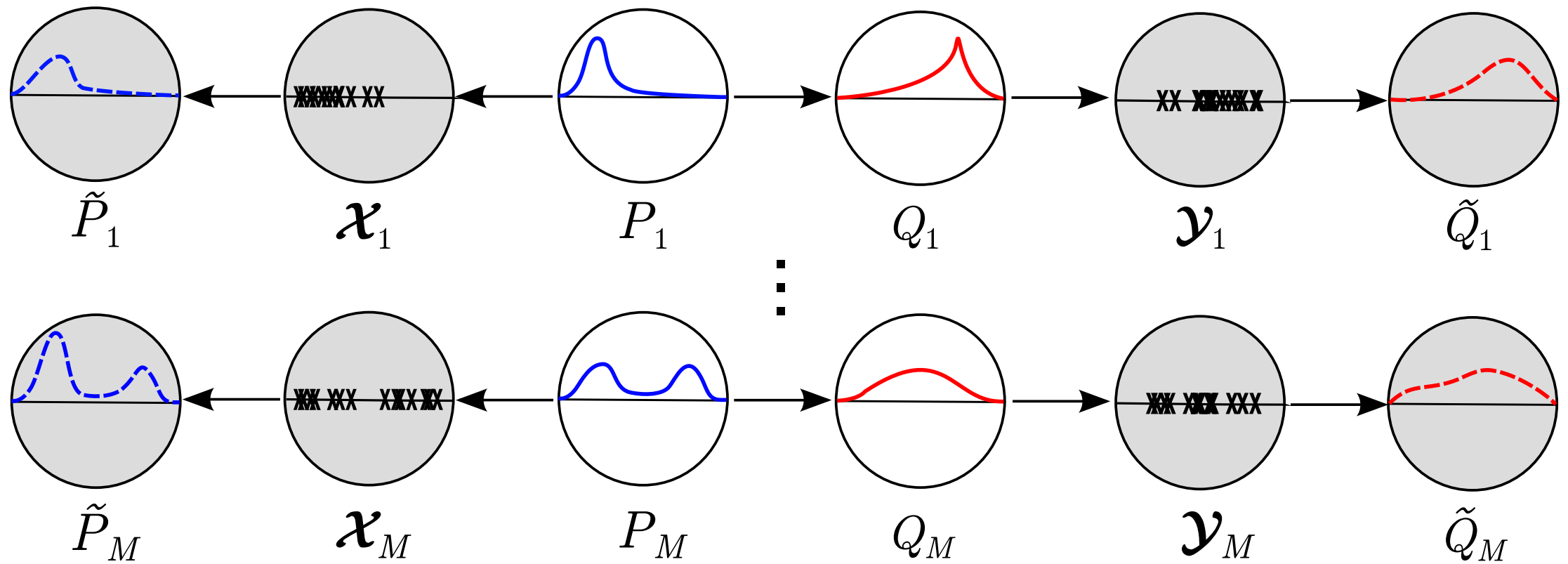
Train



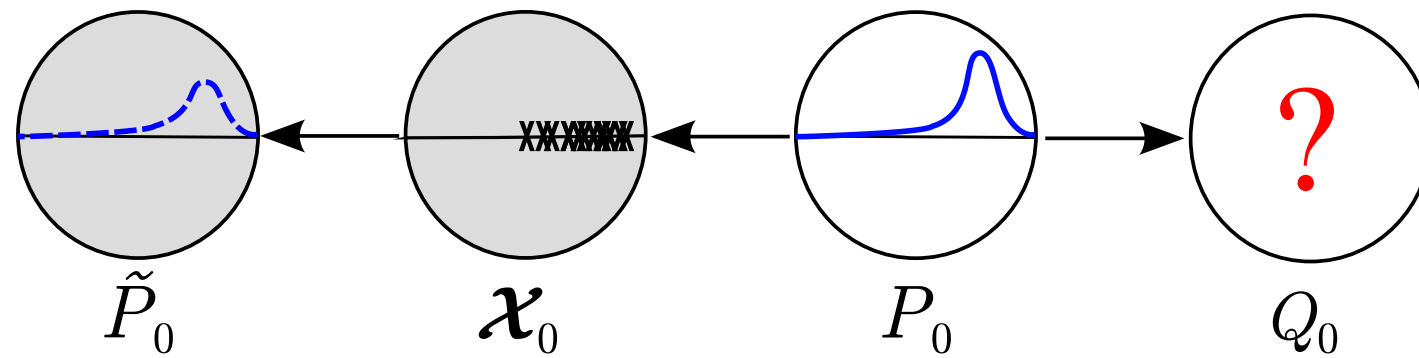
Test



Train



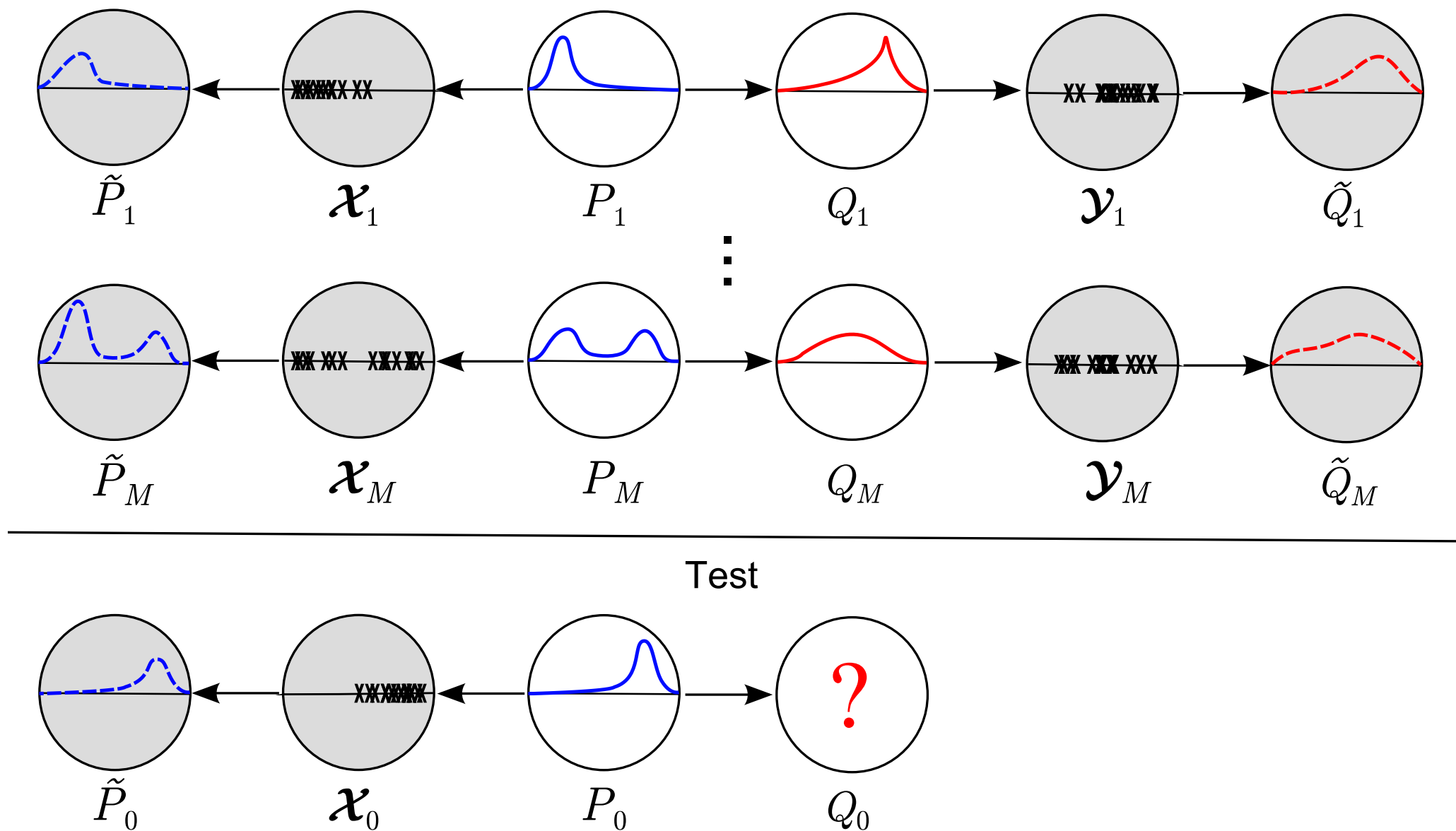
Test



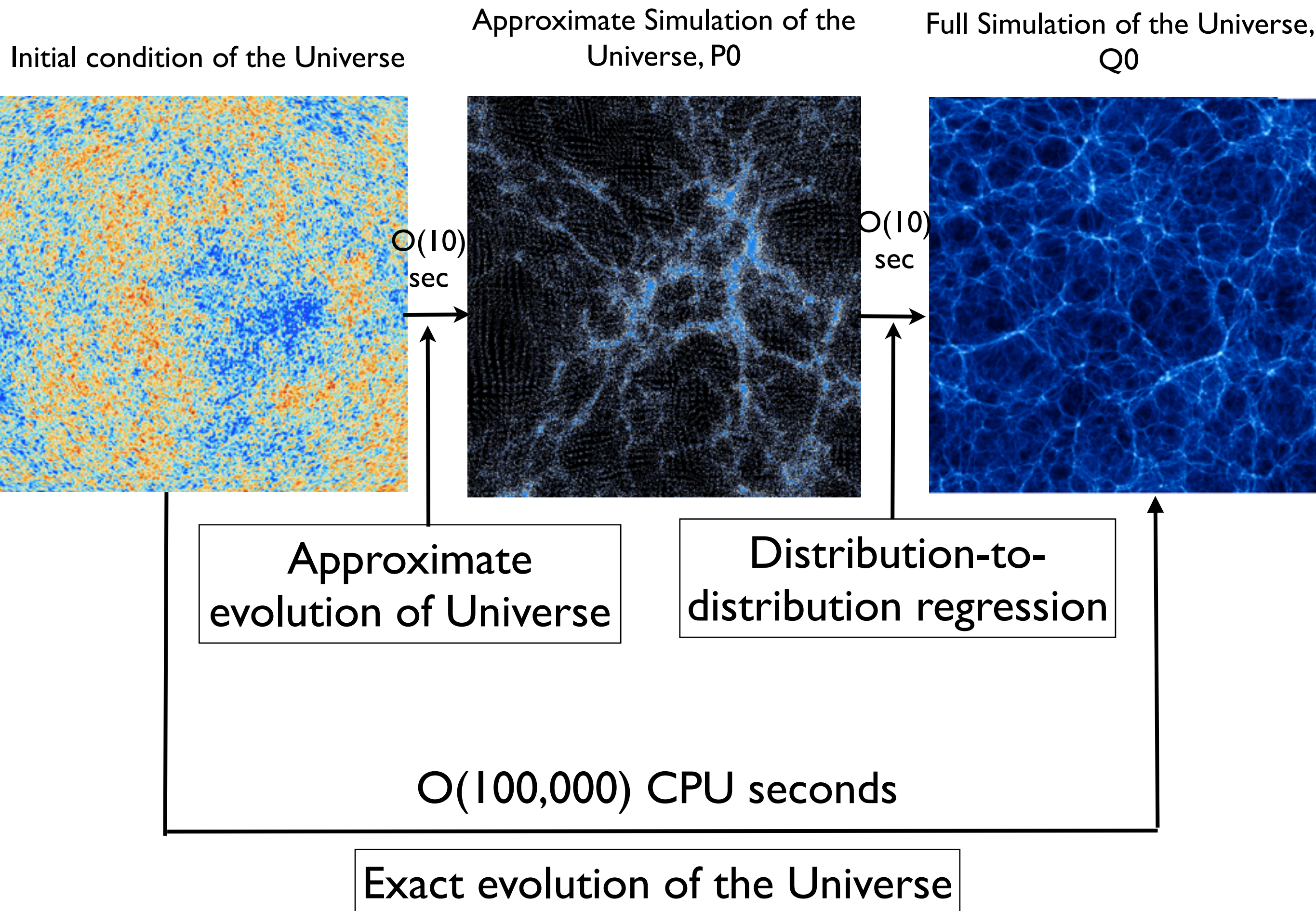
$$W(\tilde{p}_1, \tilde{p}_0) \times \tilde{q}_1 + \dots + W(\tilde{p}_M, \tilde{p}_0) \times \tilde{q}_M = \hat{f}(\tilde{p}_0)$$

(2LPT) Distribution to (N-body) Distribution Machine Learning algorithms

2LPT (P) \rightarrow N-body(Q)

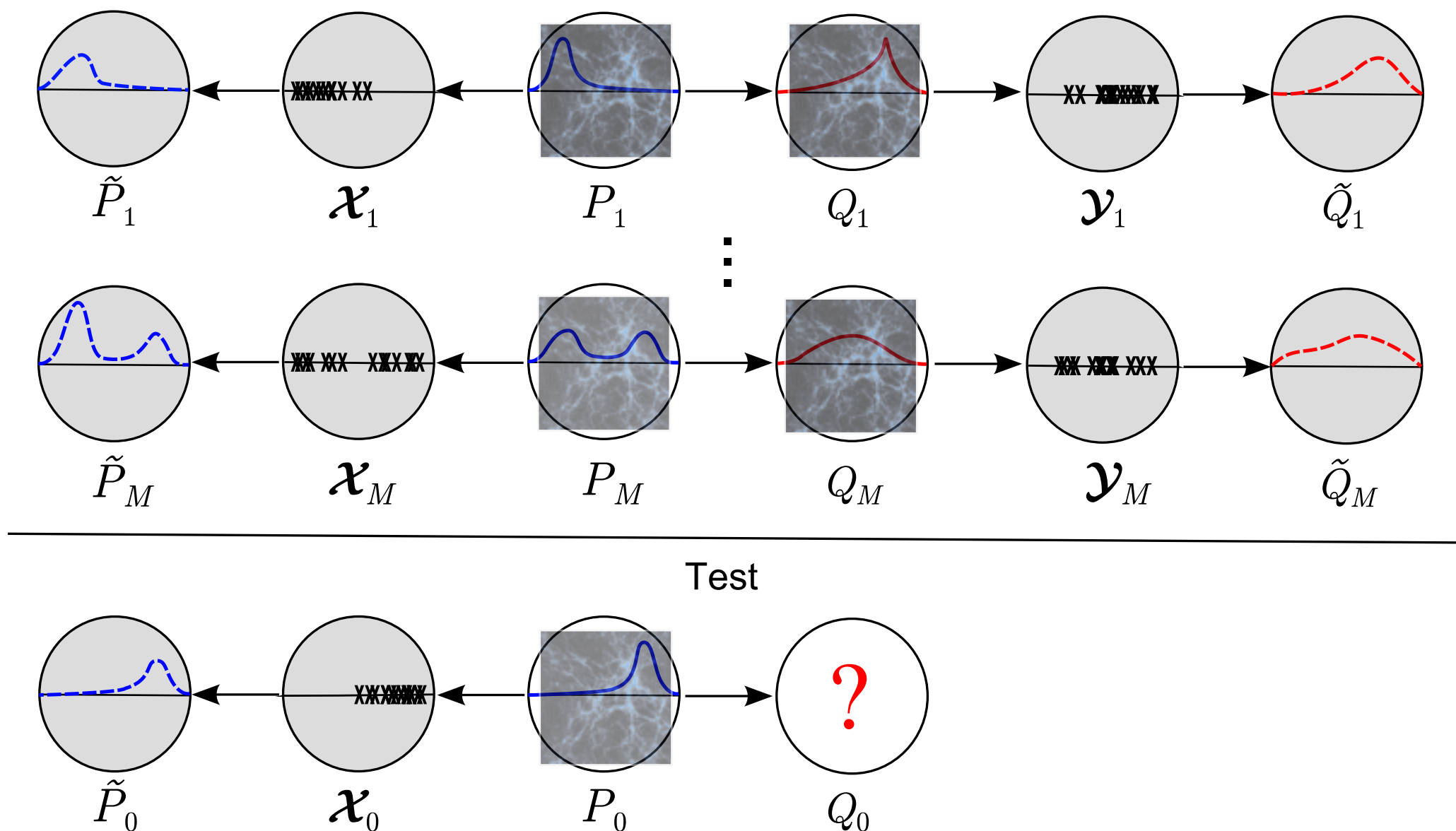


courtesy Junier Oliver



(2LPT) Distribution to (N-body) Distribution Machine Learning algorithms

2LPT (P) \rightarrow N-body(Q)



Summary

- Making many similar synthetic Universes (same cosmology) of dark matter: possibly sped up by Machine Learning by a factor of 1000!
- Instead of running very high resolution simulation or extremely complicated semi-analytical modeling, we can machine learn from either high resolution sims or actual observations to create our own halo model (of not only galaxies, but may be many other things that we have trouble modeling!)

Observations:

flux(x,y,band/wavelength), observational systematics(x,y,band/wavelength): psf, sky, dust, airmass and respective errors

Basic Reduction pipeline

Extracted object, observational systematics properties

LSS systematics removal, statistics estimation

Large Scale Structure (BAO, full clustering measurements)

Theory, Beyond Linear
Model Predictions

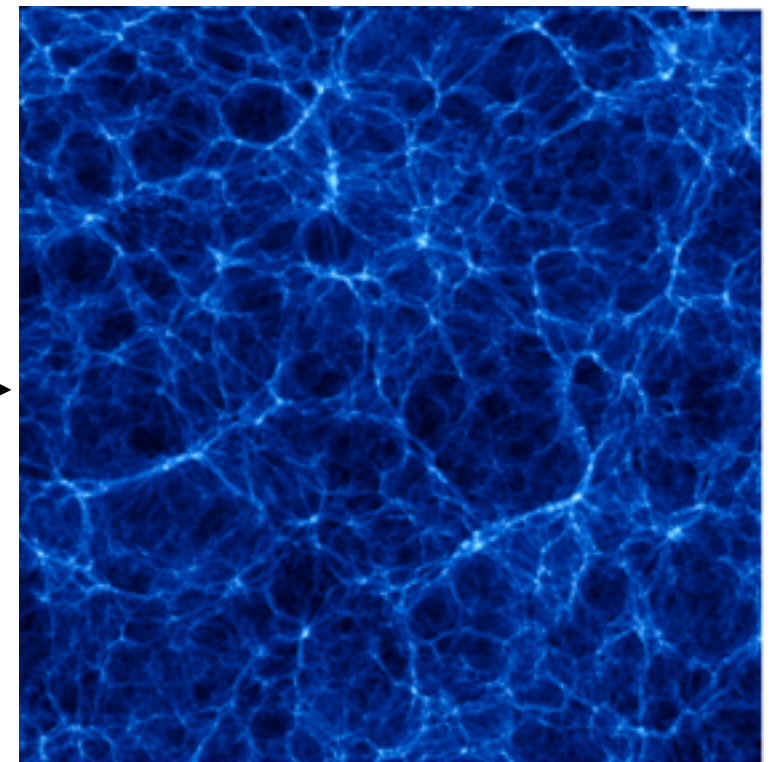
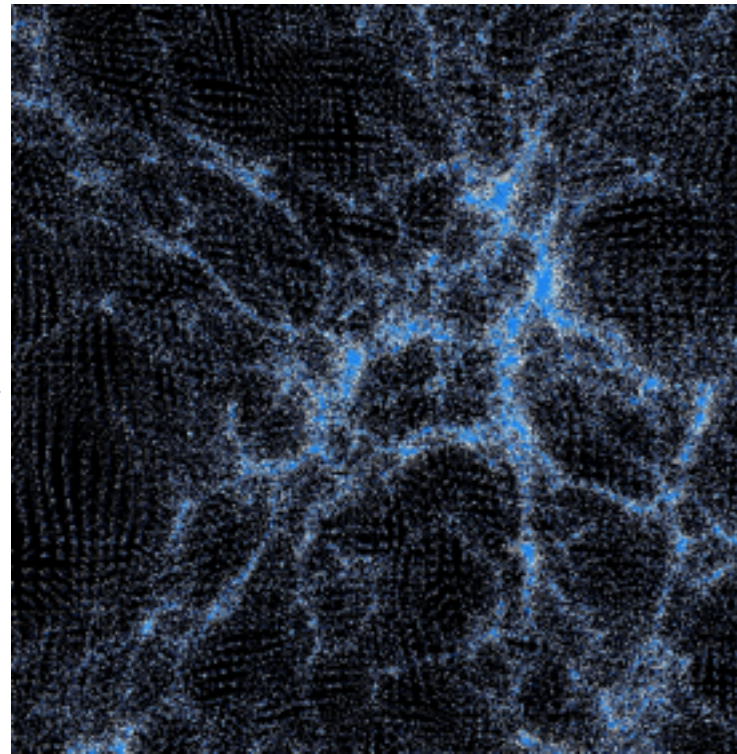
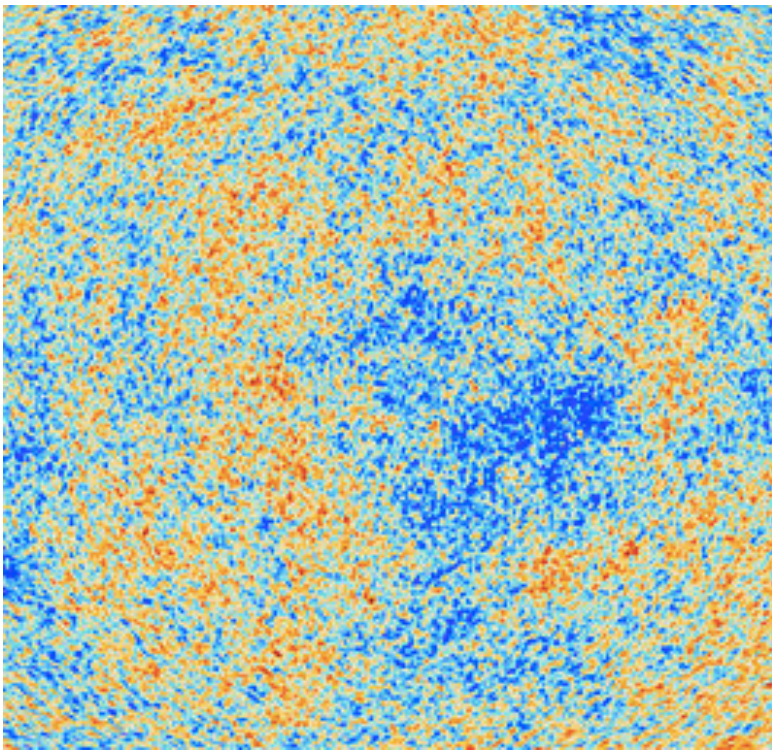
Simulations (possibly with
DM simulations + HOD),
Covariance Matrix

Cosmology (cosmological parameters, formation and
evolution of galaxies, quasars...)

Initial condition of
the Universe

Approximate
Simulation of the
Universe, P_0

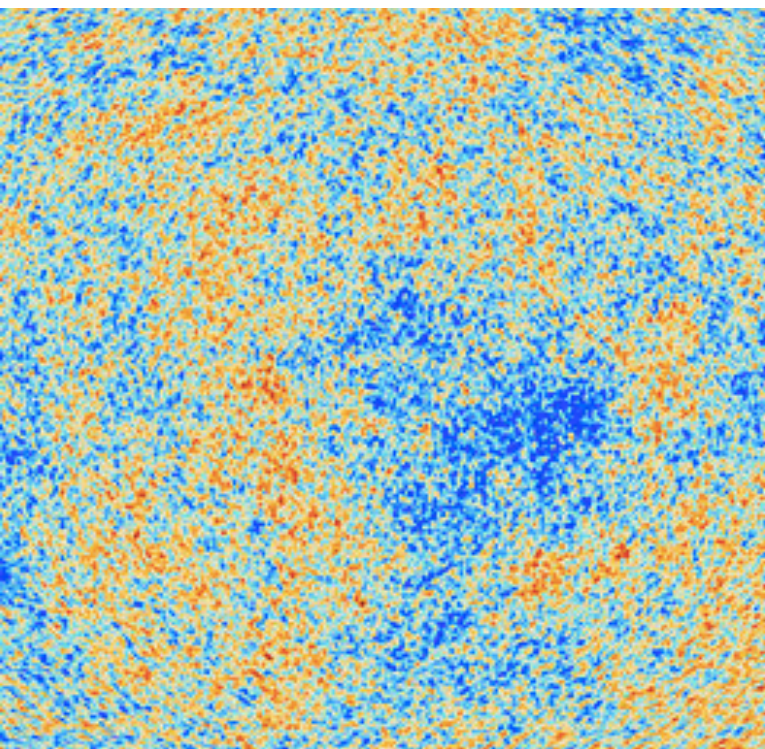
Full Simulation of
the Universe, Q_0



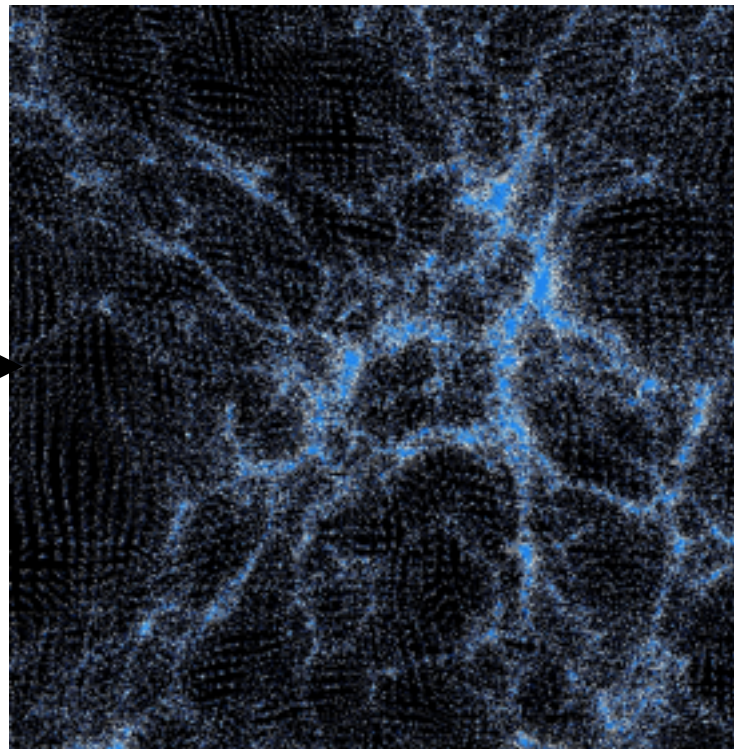
Approximate evolution
of Universe with simple
Physics

Machine Learning
Distribution-to-distribution
regression

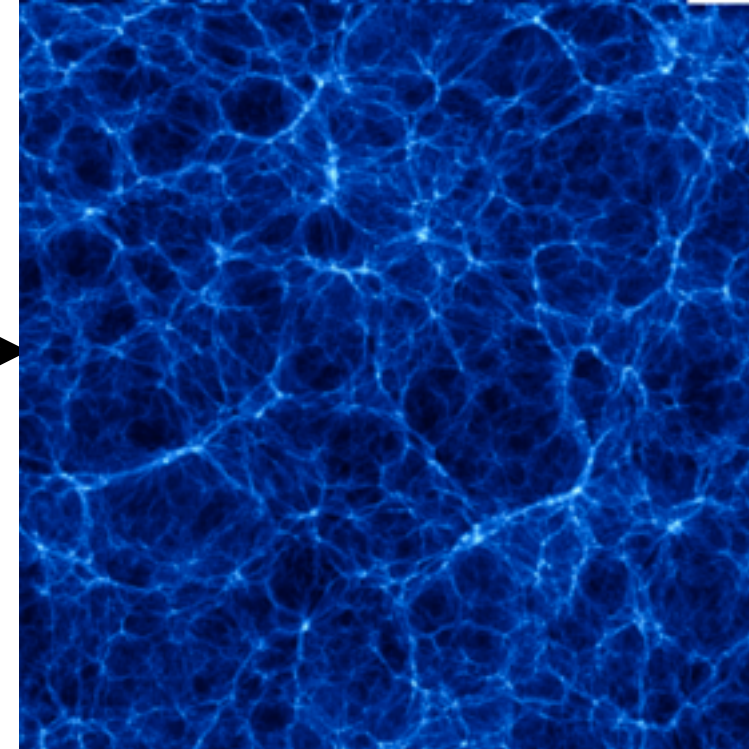
Initial condition of the Universe



Approximate Simulation of the
Universe, P0



Full Simulation of the Universe,
Q0



fast

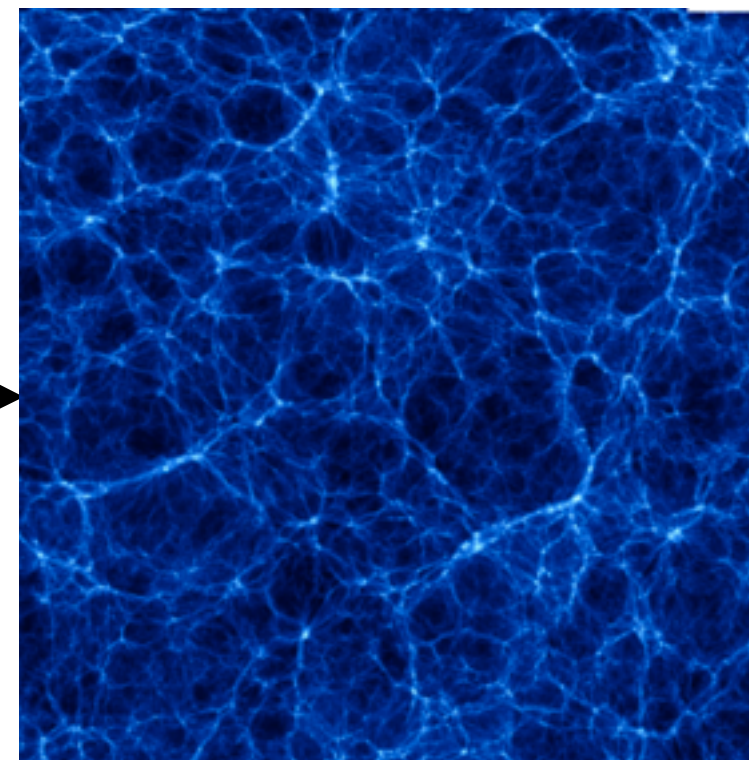
fast

Approximate
evolution of Universe

Distribution-to-distribution
regression

$O(100,000)$ CPU seconds

Exact evolution of the Universe



Initial condition of the Universe

Approximate Simulation of the Universe, P0

Full Simulation of the Universe, Q0

